

# BDMA - Massive Graph Management and Analytics

Jose Antonio Lorenzo Abril

Fall 2023



Professor: Nacéra Seghouani

Student e-mail: [jose-antonio.lorenco-abril@student-cs.fr](mailto:jose-antonio.lorenco-abril@student-cs.fr)

This is a summary of the course *Massive Graph Management and Analytics* taught at the Université Paris Saclay - CentraleSupélec by Professor Nacéra Seghouani in the academic year 23/24. Most of the content of this document is adapted from the course notes by Seghouani, [1], so I won't be citing it all the time. Other references will be provided when used.

## Contents

<b>1</b>	<b>Introduction</b>	<b>4</b>
<b>2</b>	<b>Preliminaries</b>	<b>4</b>
2.1	Graph Theory Preliminaries . . . . .	4
2.1.1	Breadth First Search (BFS) . . . . .	6
2.1.2	Depth First Search (DFS) . . . . .	9
2.1.3	Graph Representations . . . . .	10
2.1.4	Exercises . . . . .	10
2.2	Linear Algebra Preliminaries . . . . .	14
2.2.1	Exercises . . . . .	23
<b>3</b>	<b>Random Walks on Graphs</b>	<b>26</b>
3.1	First Perron-Frobenius Theorem . . . . .	26
3.2	Random Walks on Graphs . . . . .	26
3.2.1	The Stationary Distribution . . . . .	28
3.2.2	Balance Condition . . . . .	29
3.2.3	Hitting Time . . . . .	29
3.2.4	Lazy Random Walk . . . . .	30
3.3	PageRank . . . . .	30
<b>4</b>	<b>Centrality Measures</b>	<b>31</b>
4.1	Degree Centrality . . . . .	31
4.2	Neighbourhood centrality . . . . .	32
4.3	Eigenvector Centrality . . . . .	32
4.4	PageRank Centrality . . . . .	33
4.5	Katz (or alpha) centrality . . . . .	34
4.6	Clustering Coefficient Centrality . . . . .	34
4.7	Closeness Centrality . . . . .	35
4.8	Betweenness Centrality . . . . .	35

<b>5 MapReduce Computation Model</b>	<b>37</b>
5.1 The Map Function	37
5.2 Shuffling/Grouping Function	38
5.3 Reduce Function	38
5.4 MapReduce Execution Model	38
5.4.1 Coping with Node Failures	39
5.4.2 Algorithms by MapReduce	39
5.5 Use-Case: Matrix-Vector Multiplication by MapReduce	39
5.5.1 Matrix Multiplication	40
5.6 Relational Algebra by MapReduce	41
5.6.1 Selection	41
5.6.2 Projection	41
5.6.3 Join	42
5.6.4 Aggregation	42
5.7 Some Issues of MapReduce	42
<b>6 Spark Parallel Computing Framework</b>	<b>42</b>
6.1 Spark's Engine Properties and Components	43
6.2 Spark's Execution Architecture	43
6.3 Resilient Distributed Datasets (RDDs)	44
6.3.1 RDDs Operations	45
6.3.2 RDDs Actions	46
6.3.3 Caching	46

# 1 Introduction

Graph-structured data is at the heart of complex systems and plays a major role in our daily life, science and economy. Examples of this data are the cooperation between billions of individuals, or communication infrastructures with billions of cell phones, computers and satellites, the interactions between thousands of genes and metabolites within our cells, and so on.

Therefore, understanding its mathematical foundations, description, prediction, and eventually being able to control them is one of the major scientific challenges of the 21st century.

## 2 Preliminaries

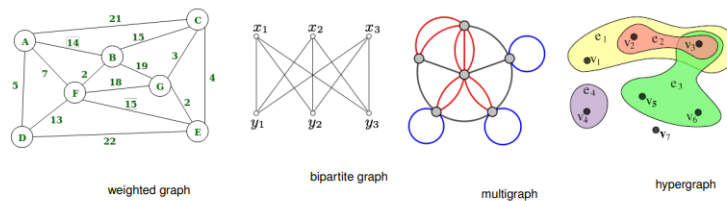
### 2.1 Graph Theory Preliminaries

A graph is a pair  $G = (V, E)$ , where  $V$  is the set of vertices and  $E \subset V \times V$  is the set of edges. Usually, we denote  $|V| = n$  and  $|E| = m$ .

There are different types of graphs:

- **Undirected:**  $(u, v) \in E \implies (v, u) \in E$ . That is, the edges goes in both directions.
- **Directed:**  $(u, v) \in E \not\implies (v, u) \in E$ . That is, the edges have direction, and it is possible that an edge goes from  $u$  to  $v$ , but not the other way.
- **Weigthed vertices:** the vertices have a weight. That is, there is a function  $w_v : V \rightarrow \mathbb{R}$ .
- **Weigthed edges:** the edges have a weight. That is, there is a function  $w_e : E \rightarrow \mathbb{R}$ .
- **Labeled vertices:** the vertices have a label,  $L_v : V \rightarrow \mathcal{L}$ , where  $\mathcal{L}$  is the set of labels.
- **Labeled edges:** the edges have a label,  $L_e : E \rightarrow \mathcal{L}$ .
- **Bipartite:** a graph  $G = (V, E)$  is bipartite if there is a partition of the vertices,  $V = V_1 \cup V_2$ , such that  $V_1 \cap V_2 = \emptyset$  and  $E = \{(v_i, v_j) | v_i \in V_1, v_j \in V_2\}$ . That is, the vertices in  $V_1$  only connect to vertices in  $V_2$ , and viceversa.
- **$k$ -Partite:** a graph  $G = (V, E)$  is  $k$ -partite if there is a  $k$ -partition of the vertices,  $V = V_1 \cup V_2 \cup \dots \cup V_k$ , such that  $V_i \cap V_j = \emptyset, \forall i \neq j$  and there is no edge  $e = (u, v)$  such that  $u, v \in V_i$ , for the same  $i$ .
- **Multigraph** or **multidigraph:** in this case, there can be several edges between two vertices. For this, we define the edges as a separate set  $E$ , and a function  $r : E \rightarrow V \times V$ , that assigns the vertices related by that edge.
- **Hypergraph:** in this case,  $E \subset 2^V$ . That is, the edges can relate 0 or more vertices. In this case, it is more appropriate to interpret  $E$  as a set of classes or hierarchies, rather than edges.
- **Complete:** a graph is complete if  $E = V \times V$ .

Some examples are:



Continuing with definitions, let  $G = (V, E)$  be a graph (directed or undirected). Let  $d_i^+$  and  $d_i^-$  denote the number of edges coming out and coming to  $v_i$ , respectively. The **degree** of  $v_i$  is

$$d_i = d_i^+ + d_i^-.$$

Note that it counts double for undirected graphs.

Now, let  $N_i^+$  and  $N_i^-$  the set of successors and predecessors of  $v_i$ , respectively. Then, the set of **neighbors** of  $v_i$  is

$$N_i = N_i^+ + N_i^-.$$

A **path** between two vertices,  $u, v \in V$ , denoted  $u \rightsquigarrow v$ , is a sequence of vertices ( $u = v_0, v_1, \dots, v_{k-1}, v_k = v$ ), where  $(v_{i-1}, v_i) \in E, \forall i = 1, \dots, k$ . The length of a path,  $L(u \rightsquigarrow v)$ , is the number of edges in the cycle, that is,  $k$ .

A **cycle** is a path from a vertex to itself,  $u \rightsquigarrow u$ .

The **distance** between two nodes,  $d(u, v)$ , is the shortest path length between them:

$$d(u, v) = \min_{u \rightsquigarrow v} L(u \rightsquigarrow v).$$

The **eccentricity** of a node,  $ecc(u)$ , is the greatest distance between  $u$  and any other vertex in the graph:

$$ecc(u) = \max_{v \in V} d(u, v).$$

Note that this could be infinity if we cannot reach some node from  $u$ . Usually, we consider only reachable nodes, because this can give us information about the graph, but a value of infinity is not very informative.

The **diameter** of a graph,  $diam(G)$ , is the greatest distance between two nodes in the graph:

$$diam(G) = \max_{u, v \in V} d(u, v) = \max_{u \in V} ecc(u).$$

The **radius** of a graph,  $rad(G)$ , is the minimum eccentricity of any vertex in the graph:

$$rad(G) = \min_{u \in V} ecc(u).$$

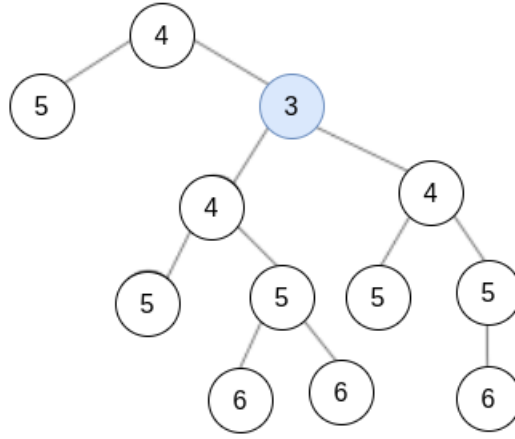
The **center** of a graph,  $C(G)$ , is the set of all vertices of minimum eccentricity, i.e., the graph radius:

$$C(G) = \{u : ecc(u) = rad(G)\}.$$

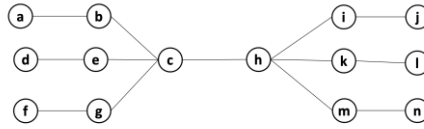
**Example 2.1.** Compute the diameter, radius and center of the following graphs:



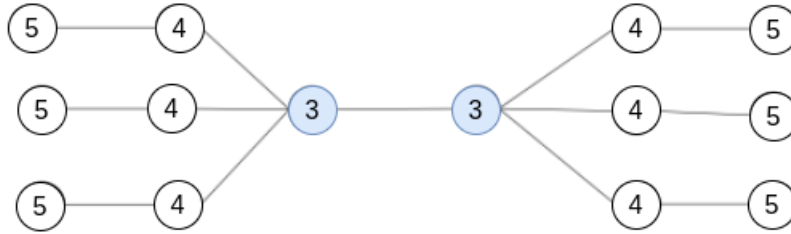
The solution is the following:



In each node, we show its eccentricity. The diameter is 6, the radius is 3 and the center is  $c$  (in blue).



Solution:



In this case, the diameter is 5, the radius is 3 and the center is  $\{c, h\}$ .

A **partial graph** of  $G = (V, E)$  is a graph  $G' = (V, E')$ , where  $E' \subset E$ .

A **subgraph** of  $G = (V, E)$  is a graph  $G' = (V', E')$  where  $V' \subset V$  and  $E' \subset E$ . Note that partial graphs are also subgraphs.

A graph  $G = (V, E)$  is said to be **connected** if, and only if,  $\forall u, v \in V, \exists u \rightsquigarrow v$ .

A (strongly) **connected component** of  $G = (V, E)$  is a subgraph  $G_{cc} = (V_{cc}, E_{cc})$ , where  $\forall u, v \in V_{cc}, \exists u \rightsquigarrow v \in V_{cc}$ . That is, a connected subgraph. It is called strongly when the paths are directed.

A graph  $G = (V, E)$  is a **tree** if, and only if,  $G$  is a connected graph without cycles. In this case, the graph has  $m = n - 1$  edges.

A graph  $G = (V, E)$  is a **forest** if, and only if, all connected components of  $G$  are trees.

### 2.1.1 Breadth First Search (BFS)

BFS is a method to traverse the nodes of a graph, by starting at one node and traversing all its neighbours. Then, all neighbours of its neighbours, and so on.

For this, we use a FIFO queue. The algorithm is:

```

1 procedure BFS(G=(V,E), r)
2   Q <- emptyset
3   enqueue(Q,r)
4   r.label = True
5
6   while Q is not empty do
7     v <- dequeue(Q)
8     for neig in neighbours(v) do
9       if not neig.label then
10        enqueue(Q,neig)
11        neig.label = True
12      end if
13    end for
14  end while
15 end procedure

```

**Example 2.2.** Apply BFS in the following graph, starting at node A.



$Q=[A]$ . We visit A's neighbours:



$Q=[F,G]$ . Now, by lexicographical order, we visit F's neighbours:



$Q=[G,B,E]$ . Now, we visit  $G$ 's neighbours. Since it has no new unvisited neighbours, there is no change.

$Q=[B,E]$ . Now, we visit  $B$ 's neighbours:



$Q=[E,C]$ . Now, we visit  $E$ 's neighbours:



$Q=[C,D]$ . Everything is visited, so the queue will be slowly emptied!



### 2.1.2 Depth First Search (DFS)

In the case of DFS, the objective is also to traverse the whole graph. The difference is that in this case we try to go as deep as we can in the graph before visiting more neighbours.

It can be implemented with a stack, let it be an explicit stack, or an implicit one.

The implementation with an explicit stack is the following:

```

1 procedure DFS(G=(V,E),r)
2   S <- emptyset
3   push(S,r)
4   while S is not empty do
5     v <- pop(S)
6     if not v.label then
7       v.label = true
8       for neig in neighbours(v) do
9         push(S, neig)
10      end for
11    end if
12  end while
13 end procedure

```

The implementation with an implicit stack is recursive, and is as follows:

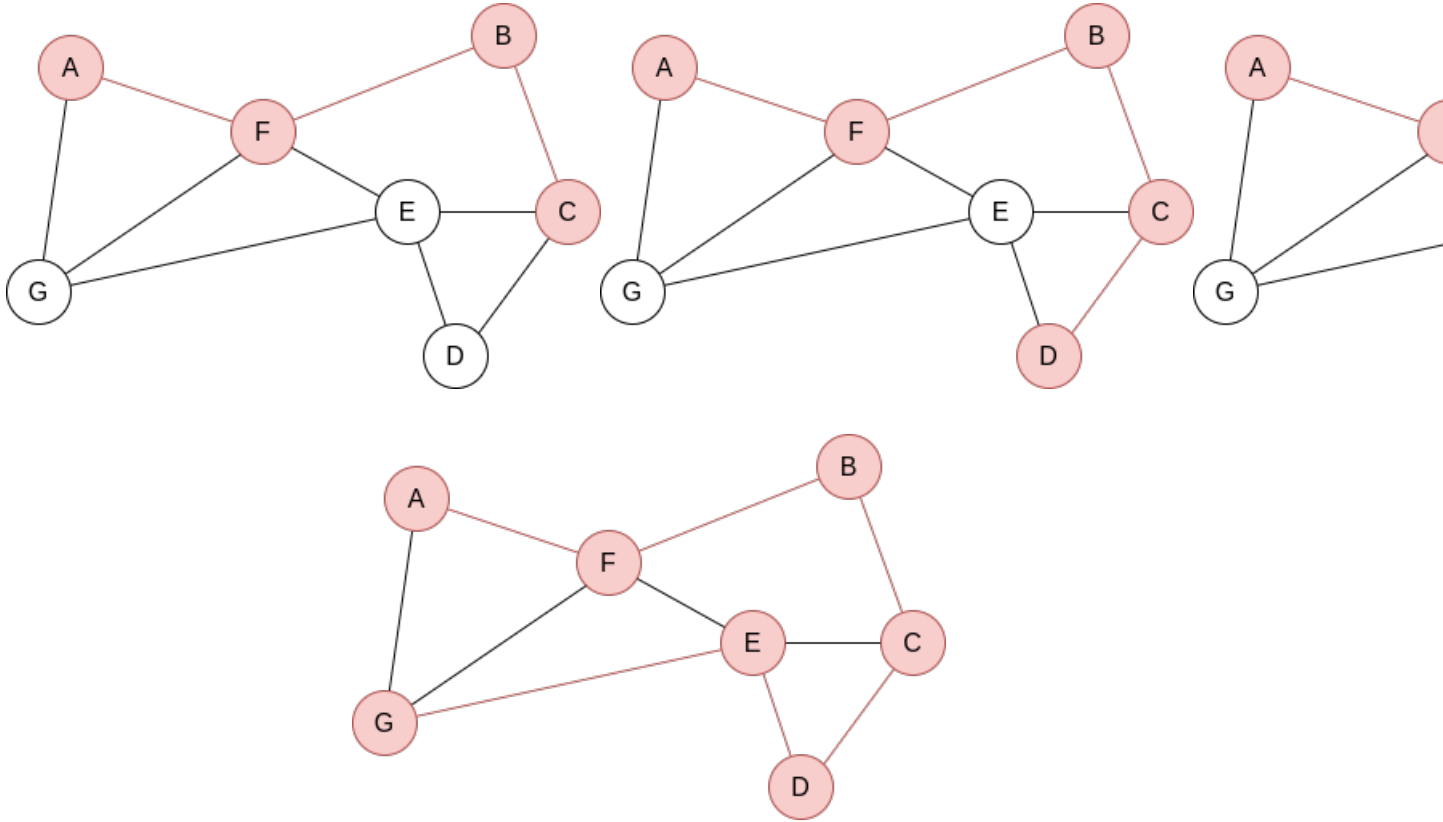
```

1 procedure BFS(G=(V,E), r)
2   Q <- emptyset
3   enqueue(Q,r)
4   r.label = True
5
6   while Q is not empty do
7     v <- dequeue(Q)
8     for neig in neighbours(v) do
9       if not neig.label then
10        enqueue(Q,neig)
11        neig.label = True
12      end if
13    end for
14  end while
15
16  whileDFS*(G=(V,E), r)
17  r.label = true
18  for neig in neighbours(r) do
19    if not neig.label then
20      DFS*(G, neig)
21    end if
22  end for
23 end procedure

```

**Example 2.3.** Let's repeat the example, now using DFS:





### 2.1.3 Graph Representations

A graph,  $G = (V, E)$ , with  $n$  vertices and  $m$  edges can be encoded using different structures:

- **Adjacency matrix:** a matrix  $A \in \mathcal{M}_{n \times n}$ , defined by

$$A_{ij} = \begin{cases} 1 & \text{if } (v_i, v_j) \in E \\ 0 & \text{otherwise} \end{cases}.$$

The adjacency matrix is symmetric for undirected graphs.

- **Adjacency list:** a list  $L$  of length  $n$  in which each vertex holds a list of its neighbours:

$$\forall u \in V, L_u = \{v \mid (u, v) \in E\}.$$

If  $G$  is directed, the choice of the direction depends on the analytic needs.

- **Incidence matrix:** a matrix  $B \in \mathcal{M}_{n \times m}$ , defined by

$$B_{ij} = \begin{cases} 1 & \text{if } e_j = (v_i, v_k) \in E \\ 0 & \text{otherwise} \end{cases}.$$

### 2.1.4 Exercises

1. Using graph traversal algorithms, propose an algorithm that computes the number of edges between a given vertex and all other vertices.

```

1 procedure n_edges(G=(V,E), r)
2   Q <- emptyset
3   enqueue(Q,r)
4   r.n_edges = 0
5
6   while Q is not empty do
7     v <- dequeue(Q)
8     for neig in neighbours(v) do
9       if not neig.n_edges then
10        enqueue(Q,neig)
11        neig.n_edges = v.n_edges + 1
12      end if
13    end for
14  end while
15 end procedure

```

2. Given the following cycles with even and odd lengths (with the distances or depths from the grey vertex), what do you think about the case of graphs with an odd cycle (in number of edges)? Is this a characteristic property? State the general case.



**Proposition:** a graph contains a cycle  $C$  with an odd number of edges if, and only if,  $\exists (x,y) \in E \mid \text{depth}(x) \neq \text{depth}(y)$ .

*Proof:* first, we know that all edges connect vertices of 'neighbouring' depths. That is,  $\forall (x,y) \in E$ , it holds  $|\text{depth}(x) - \text{depth}(y)| \leq 1$ .

[ $\Rightarrow$ ] By reduction ad absurdum, seeking a contradiction, suppose that  $\forall (x,y) \in C$ , with  $\text{depth}(x) \neq \text{depth}(y)$ . This means that  $\text{depth}(x) = \text{depth}(y) \pm 1$ . Therefore, there is, along the cycle, a node of even depth, followed by a node of odd depth, and so on. When we close the cycle, the final node is the initial one, so its depth is 0 (even). Therefore, we need an even number of edges, to conserve the parity.

[ $\Leftarrow$ ] If there is an edge  $(x,y) \in E$  with  $\text{depth}(x) \neq \text{depth}(y)$ , then we can consider the path tree that was used to annotate the depths. In this tree,  $x$  and  $y$  have a first ancestor  $z$  in common, from which we can form an odd cycle of size  $2 \cdot (\text{depth}(x) - \text{depth}(z)) + 1$  by adding the edge  $(x,y)$  to this subtree starting at  $z$ .

3. Propose an algorithm that determines if a graph contains an odd cycle.

```

1 procedure hasOddCycle(G=(V,E))
2   v <- a vertex from V
3   depths <- n_edges(G,v) #from the first exercise
4
5   for (u,v) in E do
6     if depth[u] == depth[v] then
7       return True
8     end if
9   end for
10 end procedure

```

4. In a bipartite graph, can there be a cycle with an odd number of edges? Is this a characteristic property? No, it is not possible!

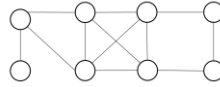
**Proposition:** A graph is bipartite if, and only if, all cycles are of even size.

[ $\Rightarrow$ ] If the graph is bipartite, any path alternates between each vertex of each partition to create a cycle ending by the initial vertex. Therefore, all cycles must be of even size.

[ $\Leftarrow$ ] Consider the partition of vertices with even depth  $V_1$ , and the partition of vertices with odd depth  $V_2$ .

Since there is no odd cycle, then, from question 2, we know that  $\forall (u,v) \in E$  it is  $\text{depth}(u) = \text{depth}(v) \pm 1$ . Therefore, the graph is bipartite.

5. Propose an algorithm that allows to determine if a graph is bipartite. Test your algorithm in the following graph. Is it bipartite? Justify your answer.



The algorithm is the same as in exercise 3, because of exercise 4.

The proposed graph is clearly not bipartite, because there are several odd cycles.

6. Graph coloring is a way of coloring the vertices of a graph in such a way that no two adjacent vertices share the same color. A 2-colorable graph is a graph that can be colored with only 2 colors.

- (a) What is the link with the previous exercise? Justify your answer.

**Proposition:** a graph is 2-colorable if, and only if, it is bipartite.

*Proof:*  $[ \implies ]$  If it is 2-colorable, with colors red and blue. Then we take  $V_1 = \{u | \text{color}(u) = \text{blue}\}$  and  $V_2 = \{u | \text{color}(u) = \text{red}\}$ .  $G$  is clearly bipartite with this partition.

$[ \impliedby ]$  If it is bipartite, with partition  $V_1$  and  $V_2$ , then we can color all nodes in  $V_1$  in blue, and all nodes in  $V_2$  in red. The graph is 2-colorable.

- (b) We want to write an algorithm, inspired by DFS search, which takes as input a graph,  $G = (V, E)$ , and which returns a pair  $(\text{result}, \text{color})$  where  $\text{result}$  is True if the graph is colorable, False otherwise, and  $\text{color}$  is a dictionary associating a color 0 or 1 to each vertex. This algorithm should stop as soon as possible when the graph is not 2-colorable.

```

1 procedure coloring(G=(V,E), r)
2   color <- {r: 0}
3   stack <- emptyset
4   push(stack, r)
5
6   while stack is not empty do
7     v <- pop(stack)
8     for neig in neighbours(v) do
9       if neig is not in color.keys then
10        push(stack, neig)
11        color[neig] = 1 - color[v]
12      elif color[neig] = color[v] then
13        return False, color
14      end if
15    end for
16  end while
17  return True, color
18 end procedure

```

7. Compute the shortest path in the following graph using Dijkstra's algorithm, starting at A:

```

1 procedure dijkstra(G=(V,E), r)
2   dist <- {r:0}
3   P <- emptyset
4
5   for v in V-{r} do
6     dist[v] = infinity
7   end for
8
9   while V-P is not empty do
10    w <- select(v in V-P and dist[v]=min_u dist[u])
11    P <- P union {w}
12    for neig in neighbours(w)-P do
13      if dist[w]+weight(neig,w) < dist[neig] then
14        dist[neig] <- dist[w]+weight(neig,w)
15      end if
16    end for
17  end while
18 end procedure

```



We start with:  $\text{dist} = \begin{array}{c|cccccc} \hline A & B & C & D & E & F & G \\ \hline 0 & \infty & \infty & \infty & \infty & \infty & \infty \\ \hline \end{array}$ , and  $P = \emptyset$ .

Now,  $w = A$  and  $P = \{A\}$ . We update  $\text{dist} = \begin{array}{c|cccccc} \hline A & B & C & D & E & F & G \\ \hline 0 & \infty & \infty & \infty & \infty & 3 & 1 \\ \hline \end{array}$ .

Now,  $w = G$  and  $P = \{A, G\}$ . We update  $\text{dist} = \begin{array}{c|cccccc} \hline A & B & C & D & E & F & G \\ \hline 0 & \infty & \infty & \infty & 3 & 2 & 1 \\ \hline \end{array}$ .

Now,  $w = F$  and  $P = \{A, F, G\}$ . We update  $\text{dist} = \begin{array}{c|cccccc} \hline A & B & C & D & E & F & G \\ \hline 0 & 6 & \infty & \infty & 3 & 2 & 1 \\ \hline \end{array}$ .

Now,  $w = E$  and  $P = \{A, E, F, G\}$ . We update  $\text{dist} = \begin{array}{c|cccccc} \hline A & B & C & D & E & F & G \\ \hline 0 & 6 & 4 & 8 & 3 & 2 & 1 \\ \hline \end{array}$ .

Now,  $w = C$  and  $P = \{A, C, E, F, G\}$ . We update  $\text{dist} = \begin{array}{c|cccccc} \hline A & B & C & D & E & F & G \\ \hline 0 & 6 & 4 & 6 & 3 & 2 & 1 \\ \hline \end{array}$ .

Now,  $w = B$  and  $P = \{A, B, C, E, F, G\}$ .  $\text{dist}$  does not change.

Finally,  $w = D$  and  $P = \{A, B, C, D, E, F, G\}$ .  $\text{dist}$  does not change.

8. Given the following graphs:



(a) Give the different representations of these graphs.

$$A_1 = \begin{array}{c} 0 \\ 1 \\ 2 \\ 3 \\ 4 \\ 5 \end{array} \begin{pmatrix} 0 & 1 & 2 & 3 & 4 & 5 \\ 0 & 1 & 0 & 0 & 1 & 1 \\ 1 & 0 & 1 & 1 & 0 & 1 \\ 0 & 1 & 0 & 1 & 0 & 1 \\ 0 & 1 & 1 & 0 & 0 & 1 \\ 1 & 0 & 0 & 0 & 0 & 1 \\ 1 & 1 & 1 & 1 & 1 & 0 \end{pmatrix}, L_1 = \begin{array}{l} 0 : \{1, 4, 5\} \\ 1 : \{0, 2, 3, 5\} \\ 2 : \{1, 3, 5\} \\ 3 : \{1, 2, 5\} \\ 4 : \{0, 5\} \\ 5 : \{0, 1, 2, 3, 4\} \end{array}$$

$$B_1 = \begin{array}{c} 0 \\ 1 \\ 2 \\ 3 \\ 4 \\ 5 \end{array} \begin{pmatrix} 01 & 04 & 05 & 12 & 13 & 15 & 23 & 25 & 35 & 45 \\ 1 & 1 & 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 1 & 0 & 0 & 1 & 1 & 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 & 0 & 1 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & 0 & 1 & 0 & 1 & 0 \\ 0 & 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 \\ 0 & 0 & 1 & 0 & 0 & 1 & 0 & 1 & 1 & 1 \end{pmatrix}$$

$$\begin{aligned}
A_2 = & \begin{matrix} & \begin{matrix} 0 & 1 & 2 & 3 & 4 & 5 \end{matrix} \\ \begin{matrix} 0 \\ 1 \\ 2 \\ 3 \\ 4 \\ 5 \end{matrix} & \begin{pmatrix} 0 & 1 & 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 & 0 & 1 \\ 0 & 0 & 0 & 1 & 0 & 1 \\ 0 & 0 & 0 & 1 & 0 & 0 \\ 1 & 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & 0 \end{pmatrix} \end{matrix}, L_2 = \begin{matrix} 0 : \{1, 4\} \\ 1 : \{5\} \\ 2 : \{3, 5\} \\ 3 : \{3\} \\ 4 : \{0, 1\} \\ 5 : \{4\} \end{matrix}, \\
B_1 = & \begin{matrix} & \begin{matrix} 01 & 04 & 15 & 23 & 25 & 33 & 40 & 41 & 54 \end{matrix} \\ \begin{matrix} 0 \\ 1 \\ 2 \\ 3 \\ 4 \\ 5 \end{matrix} & \begin{pmatrix} 1 & 1 & 0 & 0 & 0 & 0 & 2 & 0 & 0 \\ 2 & 0 & 1 & 0 & 0 & 0 & 0 & 2 & 0 \\ 0 & 0 & 0 & 1 & 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 2 & 0 & 3 & 0 & 0 & 0 \\ 0 & 2 & 0 & 0 & 0 & 0 & 1 & 1 & 2 \\ 0 & 0 & 2 & 0 & 2 & 0 & 0 & 0 & 1 \end{pmatrix} \end{matrix}
\end{aligned}$$

(b) Compute  $A^2, A^3$ . What does  $A_{ij}^r$  represents?

$A_{ij}^r$  represents the number of paths of length  $r$  from node  $i$  to node  $j$ .

(c) What is the complexity of  $A^r$ ? Is it possible to reduce it?

Computing  $A^r$  is  $O(rn^3)$ , since it requires  $r$  products of complexity  $O(n^3)$ .

However, we can reuse some results to reduce the complexity:

- If  $r$  is even, we can do  $A^r = (A^{\frac{r}{2}})^2$ .
- If  $r$  is odd, we can do  $A^r = A \left( A^{\frac{r-1}{2}} \right)^2$ .

Therefore, we can obtain  $A^r$  in  $O(\log r \cdot n^3)$ .

## 2.2 Linear Algebra Preliminaries

A **norm** is a function  $f$  that measures the size of a vector. It must satisfy the following properties:

- $f(x) = 0 \iff x = 0$ .
- Linear on scale factors:

$$f(\alpha x) = |\alpha| f(x), \forall \alpha \in \mathbb{R}.$$

- Triangle inequality:

$$f(x + y) \leq f(x) + f(y).$$

A widely use family of norms are the  $p$ -norms:

$$\|x\|_p = \sqrt[p]{\sum_i |x_i|^p},$$

with the most common one being the Euclidean norm, for  $p = 2$ :

$$\|x\| = \sqrt{\sum_i x_i^2}.$$

The **determinant** of a square matrix is equal to the hypervolume of the parallelotope defined by the vectors of the matrix. It is 0 if, and only if, the set of vectors is colinear.

The determinant can be used for many things:

- We can represents linear systems with matrices as  $Y = AX$ , and there are many methods to solve this efficiently.

- With the determinant we can compute the **characteristic polynomial** of  $A$ , whose roots are the eigenvalues of  $A$ .

Some properties of the determinant are:

- $|I| = 1$ , where  $I$  is the identity matrix.
- $|A| = 0$  if  $A$  is singular (not invertible).
- $|AB| = |A| |B|$ .
- $|A^T| = |A|$ .
- $|cA| = c^n |A|$ , where  $n$  is the dimension of  $A$ .

A square matrix,  $A$ , is **invertible** (non-singular, non-degenerate), with inverse denoted  $A^{-1}$ , if  $\exists B$  such that

$$AB = BA = I,$$

in this case,  $A^{-1} = B$ .

**Proposition 2.1.** *For a square matrix,  $A$ , the following properties are equivalent:*

- $A$  is invertible.
- All vectors in  $A$  are linearly independent.
- $|A| \neq 0$ .
- $A^T$  is invertible.
- $0$  is not an eigenvalue of  $A$ .

Properties of the inverse:

- $(A^{-1})^{-1} = A$ .
- $(A^T)^{-1} = (A^{-1})^T$ .
- $(AB)^{-1} = B^{-1}A^{-1}$ .
- $(cA)^{-1} = \frac{1}{c}A^{-1}$  for  $c \neq 0$ .
- $|A^{-1}| = \frac{1}{|A|}$ .

An **eigenvector** or characteristic vector of a linear transformation,  $T$ , is a non-zero vector that changes by a scalar factor,  $\lambda$ , when transformed by  $T$ . That is,  $v$  is an eigenvector of the linear transformation  $T$  if

$$T(v) = \lambda v.$$

There is a direct correspondence between  $n \times n$  matrices and linear transformation in the  $n$ -dimensional vector space into itself. That is, every linear transformation  $T$  can be represented as a matrix  $A_T$  (the matrix depends on the chosen base). Therefore, we can say that  $A_T$  has an eigenvector  $v$  if

$$A_T v = \lambda v.$$

The scale factors of the eigenvectors are called **eigenvalues**.

We can find the eigenvalues by solving a polynomial function on  $\lambda$  called the **characteristic polynomial** of  $A_T$ :

$$(A - \lambda I) v = 0.$$

Now, this equation has non-zero solution if, and only if,

$$|A - \lambda I| = 0.$$

Therefore, we can compute  $|A - \lambda I|$  and find all values of  $\lambda$  that makes it equal to 0.

Once we have the eigenvalues, we can use them to find the corresponding eigenvectors.

**Example 2.4.** Compute the eigenvalues and eigenvectors of  $A = \begin{pmatrix} 2 & 1 \\ 1 & 2 \end{pmatrix}$ .

$$|A - \lambda I| = \begin{vmatrix} 2 - \lambda & 1 \\ 1 & 2 - \lambda \end{vmatrix} = (2 - \lambda)^2 - 1 = \lambda^2 - 4\lambda + 3.$$

This has as solutions

$$\lambda = \frac{4 \pm \sqrt{16 - 12}}{2} = \frac{4 \pm 2}{2} = 2 \pm 1.$$

Therefore, we have  $\lambda_1 = 1, \lambda_2 = 3$ .

To find the eigenvectors, we solve

$$Av = \lambda v \iff \begin{cases} 2x + y = \lambda x \\ x + 2y = \lambda y \end{cases}.$$

For  $\lambda_1 = 1$ , it is

$$\begin{cases} 2x + y = x \\ x + 2y = y \end{cases} \iff x = -y,$$

so the eigenvector associated to  $\lambda_1 = 1$  is

$$v_{\lambda_1} = \begin{pmatrix} t \\ -t \end{pmatrix}.$$

For  $\lambda_2 = 3$ , it is

$$\begin{cases} 2x + y = 3x \\ x + 2y = 3y \end{cases} \iff x = y,$$

so the eigenvector associated to  $\lambda_2 = 3$  is

$$v_{\lambda_2} = \begin{pmatrix} t \\ t \end{pmatrix}.$$

We call the **algebraic multiplicity**,  $t_i$ , of the eigenvalue  $\lambda_i$  to its multiplicity as root of the characteristic polynomial:

$$P(A) = |A - \lambda I| = (\lambda - \lambda_1)^{t_1} (\lambda - \lambda_2)^{t_2} \cdot \dots \cdot (\lambda - \lambda_k)^{t_k}.$$

Note that  $A$  can have at most  $n$  distinct eigenvalues, although some of them may be complex.

**Proposition 2.2.** *If the eigenvalues of  $A$  are all different, then the corresponding eigenvectors are linearly independent.*

The **eigenspace** of an eigenvalue,  $\lambda$ , is the space generated by the eigenvectors associated to  $\lambda$ .

The dimension of the eigenspace of  $\lambda$  is the **geometric multiplicity** of  $\lambda$ . The geometric multiplicity of an eigenvalue is, at most, its algebraic multiplicity.

**Example 2.5.** Let's get some eigenspaces:



$$A = \begin{pmatrix} -1 & 1 & 0 \\ -4 & 3 & 0 \\ 1 & 0 & 2 \end{pmatrix}, \text{ so}$$

$$\begin{aligned} |A - \lambda I| &= \begin{vmatrix} -1-\lambda & 1 & 0 \\ -4 & 3-\lambda & 0 \\ 1 & 0 & 2-\lambda \end{vmatrix} = (-1-\lambda)(3-\lambda)(2-\lambda) + 4(2-\lambda) \\ &= (2-\lambda)[(-1-\lambda)(3-\lambda) + 4] = (2-\lambda)(-3 + \lambda - 3\lambda + \lambda^2 + 4) \\ &= (2-\lambda)(\lambda^2 - 2\lambda + 1) \\ &= (2-\lambda)(\lambda - 1)^2. \end{aligned}$$

This has roots  $\lambda_1 = 1$ , with algebraic multiplicity 2, and  $\lambda_2 = 2$ , with algebraic multiplicity 1.

Now, we get the eigenvectors associated to them:

$$Av = \lambda v \iff \begin{cases} -x + y &= \lambda x \\ -4x + 3y &= \lambda y \\ x + 2z &= \lambda z \end{cases}$$

For  $\lambda_1$  this is

$$\begin{cases} -x + y &= x \\ -4x + 3y &= y \\ x + 2z &= z \end{cases} \iff \begin{cases} y &= 2x \\ -4x + 3y &= y \\ x &= -z \end{cases},$$

$$\text{so } v_{\lambda_1} = \begin{pmatrix} t \\ 2t \\ -t \end{pmatrix}, \text{ with dimension 1 (it could be 2).}$$

For  $\lambda_2$  this is

$$\begin{cases} -x + y &= 2x \\ -4x + 3y &= 2y \\ x + 2z &= 2z \end{cases} \iff \begin{cases} y &= 3x \\ -4x &= -y \\ x + 2z &= \lambda z \end{cases} \iff \begin{cases} x = 0 \\ y = 0 \\ 2z = 2z \end{cases},$$

$$\text{so } v_{\lambda_2} = \begin{pmatrix} 0 \\ 0 \\ t \end{pmatrix}, \text{ with dimension 1 (it could not be differently).}$$

$$B = \begin{pmatrix} 4 & 6 & 0 \\ -3 & -5 & 0 \\ -3 & -6 & 1 \end{pmatrix}, \text{ so}$$

$$\begin{aligned} |B - \lambda I| &= \begin{vmatrix} 4-\lambda & 6 & 0 \\ -3 & -5-\lambda & 0 \\ -3 & -6 & 1-\lambda \end{vmatrix} = (4-\lambda)(-5-\lambda)(1-\lambda) + 18(1-\lambda) \\ &= (1-\lambda)[(4-\lambda)(-5-\lambda) + 18] = (1-\lambda)(-20 - 4\lambda + 5\lambda + \lambda^2 + 18) \\ &= (1-\lambda)(\lambda^2 + \lambda - 2) = (1-\lambda)^2(-2-\lambda). \end{aligned}$$

This has roots  $\lambda_1 = 1$ , with algebraic multiplicity 2, and  $\lambda_2 = -2$ , with algebraic multiplicity 1.

Now, we get the eigenvectors associated to them:

$$Av = \lambda v \iff \begin{cases} 4x + 6y &= \lambda x \\ -3x - 5y &= \lambda y \\ -3x - 6y + z &= \lambda z \end{cases}$$

For  $\lambda_1 = 1$ , we have

$$\begin{cases} 4x + 6y &= x \\ -3x - 5y &= y \\ -3x - 6y + z &= z \end{cases} \iff \begin{cases} x &= -2y \\ z &= z \end{cases}.$$

Therefore, the eigenspace associated to  $\lambda_1$  is

$$E(\lambda_1) = \left\{ \begin{pmatrix} -2t \\ t \\ 0 \end{pmatrix}, \begin{pmatrix} 0 \\ 0 \\ u \end{pmatrix} \right\}.$$

For  $\lambda_2 = 2$ , we have

$$\begin{cases} 4x + 6y &= -2x \\ -3x - 5y &= -2y \\ -3x - 6y + z &= -2z \end{cases} \iff \begin{cases} x &= -y \\ -3y + z &= -2z \end{cases} \iff \begin{cases} x &= -y \\ y &= z \end{cases}.$$

Thus, the eigenspace associated to  $\lambda_2$  is

$$E(\lambda_2) = \begin{pmatrix} -t \\ t \\ t \end{pmatrix}.$$

$$C = \begin{pmatrix} 1 & -1 & 0 \\ -1 & 2 & 1 \\ 0 & 1 & 1 \end{pmatrix},$$

$$\begin{aligned} |C - \lambda I| &= \begin{vmatrix} 1 - \lambda & -1 & 0 \\ -1 & 2 - \lambda & 1 \\ 0 & 1 & 1 - \lambda \end{vmatrix} = (1 - \lambda)^2 (2 - \lambda) - 2(1 - \lambda) \\ &= (1 - \lambda) [(1 - \lambda)(2 - \lambda) - 2] \\ &= (1 - \lambda) (2 - \lambda - 2\lambda + \lambda^2 - 2) \\ &= (1 - \lambda) (\lambda^2 - 3\lambda) \\ &= (1 - \lambda) (\lambda - 3) \lambda. \end{aligned}$$

$$Cv = \lambda v \iff \begin{cases} x - y &= \lambda x \\ -x + 2y + z &= \lambda y \\ y + z &= \lambda z \end{cases}.$$

$\lambda_1 = 0$ :

$$\begin{cases} x - y &= 0 \\ -x + 2y + z &= 0 \\ y + z &= 0 \end{cases} \iff \begin{cases} x = y \\ y + z &= 0 \end{cases} \iff \begin{cases} x = y \\ y = -z \end{cases},$$

so

$$E(\lambda_1) = \begin{pmatrix} t \\ t \\ -t \end{pmatrix}.$$

$\lambda_2 = 1$ :

$$\begin{cases} x - y &= x \\ -x + 2y + z &= y \\ y + z &= z \end{cases} \iff \begin{cases} y = 0 \\ -x + z &= \\ z &= z \end{cases} \iff \begin{cases} y = 0 \\ x = z \end{cases},$$

so

$$E(\lambda_2) = \begin{pmatrix} t \\ 0 \\ t \end{pmatrix}.$$

$\lambda_3 = 3$ :

$$\begin{cases} x - y &= 3x \\ -x + 2y + z &= 3y \\ y + z &= 3z \end{cases} \iff \begin{cases} y &= -2x \\ y &= 2z \end{cases},$$

so

$$E(\lambda_3) = \begin{pmatrix} -t \\ 2t \\ t \end{pmatrix}.$$

$$D = \begin{pmatrix} 1 & -1 & 4 \\ 3 & 2 & -1 \\ 2 & 1 & -1 \end{pmatrix},$$

$$\begin{aligned} |D - \lambda I| &= \begin{vmatrix} 1 - \lambda & -1 & 4 \\ 3 & 2 - \lambda & -1 \\ 2 & 1 & -1 - \lambda \end{vmatrix} \\ &= (1 - \lambda)(2 - \lambda)(-1 - \lambda) + 12 + 2 - 8(2 - \lambda) + 1 - \lambda + 3(-1 - \lambda) \\ &= (2 - 3\lambda + \lambda^2)(-1 - \lambda) - 4 + 4\lambda \\ &= -2 - 2\lambda + 3\lambda + 3\lambda^2 - \lambda^2 - \lambda^3 - 4 + 4\lambda \\ &= -\lambda^3 + 2\lambda^2 + 5\lambda - 6 \end{aligned}$$

To obtain the roots, we can use Ruffini:

$$\begin{array}{r|rrrr} & -1 & 2 & 5 & -6 \\ 1 & & -1 & 1 & 6 \\ \hline & -1 & 1 & 6 & 0 \end{array}$$

So  $\lambda_1 = 1$  is a root and we have now  $-\lambda^2 + \lambda + 6 = 0$ , obtaining

$$\lambda = \frac{-1 \pm \sqrt{1 + 24}}{-2} = \frac{-1 \pm 5}{-2},$$

and we get  $\lambda_2 = -2$  and  $\lambda_3 = 3$ .

$$Dv = \lambda v \iff \begin{cases} x - y + 4z &= \lambda x \\ 3x + 2y - z &= \lambda y \\ 2x + y - z &= \lambda z \end{cases}.$$

$\lambda_1 = 1$ :

$$\begin{cases} x - y + 4z &= x \\ 3x + 2y - z &= y \\ 2x + y - z &= z \end{cases} \iff \begin{cases} y &= 4z \\ 3x + 3z &= 0 \\ 2x + 2z &= 0 \end{cases} \iff \begin{cases} y = 4z \\ x = -z \end{cases}.$$

$$\text{Then, } E(\lambda_1) = \begin{pmatrix} -t \\ 4t \\ t \end{pmatrix}.$$

$\lambda_2 = -2$ :

$$\begin{aligned} \begin{cases} x - y + 4z &= -2x \\ 3x + 2y - z &= -2y \\ 2x + y - z &= -2z \end{cases} &\iff \begin{cases} 3x - y + 4z &= 0 \\ 3x + 4y - z &= 0 \\ 2x + y + z &= 0 \end{cases} \iff \begin{cases} 5y - 5z &= 0 \\ 2x + y + z &= 0 \end{cases} \\ &\iff \begin{cases} y = z \\ 2x + 2y &= 0 \end{cases} \iff \begin{cases} y = z \\ x = -y \end{cases}. \end{aligned}$$

Then,  $E(\lambda_2) = \begin{pmatrix} -t \\ t \\ t \end{pmatrix}$ .

$\lambda_3 = 3$ :

$$\begin{cases} x - y + 4z = 3x \\ 3x + 2y - z = 3y \\ 2x + y - z = 3z \end{cases} \iff \begin{cases} -2x - y + 4z = 0 \\ 3x - y - z = 0 \\ 2x + y - 4z = 0 \end{cases} \iff \begin{cases} -2x - y + 4z = 0 \\ 5x - 5z = 0 \end{cases} \\ \iff \begin{cases} y = 2z \\ x = z \end{cases}.$$

Then,  $E(\lambda_3) = \begin{pmatrix} t \\ 2t \\ t \end{pmatrix}$ .

$$E = \begin{pmatrix} 6 & -2 & 2 \\ -2 & 3 & -1 \\ 2 & -1 & 3 \end{pmatrix},$$

$$\begin{aligned} |E - \lambda I| &= \begin{vmatrix} 6 - \lambda & -2 & 2 \\ -2 & 3 - \lambda & -1 \\ 2 & -1 & 3 - \lambda \end{vmatrix} \\ &= (6 - \lambda)(3 - \lambda)^2 + 4 + 4 - 4(3 - \lambda) - (6 - \lambda) - 4(3 - \lambda) \\ &= (6 - \lambda)(9 - 6\lambda + \lambda^2) + 2 - 8(3 - \lambda) + \lambda \\ &= 54 - 36\lambda + 6\lambda^2 - 9\lambda + 6\lambda^2 - \lambda^3 - 22 + 8\lambda + \lambda \\ &= -\lambda^3 + 12\lambda^2 - 36\lambda + 32. \end{aligned}$$

Again, we can use the Ruffini rule:

$$\begin{array}{r|rrrr} & -1 & 12 & -36 & 32 \\ 2 & & -2 & 20 & -32 \\ \hline & -1 & 10 & -16 & 0 \end{array}$$

So  $\lambda_1 = 2$  is a root, and we now have  $-\lambda^2 + 10\lambda - 16 = 0$ , which gives us

$$\lambda = \frac{-10 \pm \sqrt{100 - 64}}{-2} = \frac{-10 \pm 6}{-2} = 5 \pm 3.$$

Therefore,  $\lambda_1$  is a double root and the other root is  $\lambda_2 = 8$ .

$$Ev = \lambda v \iff \begin{cases} 6x - 2y + 2z = \lambda x \\ -2x + 3y - z = \lambda y \\ 2x - y + 3z = \lambda z \end{cases}.$$

$\lambda_1 = 2$ :

$$\begin{cases} 6x - 2y + 2z = 2x \\ -2x + 3y - z = 2y \\ 2x - y + 3z = 2z \end{cases} \iff \begin{cases} 4x - 2y + 2z = 0 \\ -2x + y - z = 0 \\ 2x - y + z = 0 \end{cases} \iff \begin{cases} 4x - 2y + 2z = 0 \\ 2x - y + z = 0 \end{cases} \\ \iff 2x - y + z = 0$$

If  $x = 0$ :  $y = z$ .

If  $x = t \neq 0$ :  $-y + z = -2t$ , working for  $y = t$  and  $z = -t$ .

So

$$E(\lambda_1) = \left\{ \begin{pmatrix} 0 \\ t \\ t \end{pmatrix}, \begin{pmatrix} t \\ t \\ -t \end{pmatrix} \right\}.$$

$\lambda_2 = 8$ :

$$\begin{aligned} \begin{cases} 6x - 2y + 2z &= 8x \\ -2x + 3y - z &= 8y \\ 2x - y + 3z &= 8z \end{cases} &\iff \begin{cases} -2x - 2y + 2z &= 0 \\ -2x - 5y - z &= 0 \\ 2x - y - 5z &= 0 \end{cases} \iff \begin{cases} -3y - 3z &= 0 \\ 2x - y - 5z &= 0 \end{cases} \\ &\iff \begin{cases} y = -z \\ 2x - 4z &= 0 \end{cases} \iff \begin{cases} y = -z \\ x = 2z \end{cases}. \end{aligned}$$

Therefore,

$$E(\lambda_2) = \begin{pmatrix} 2t \\ -t \\ t \end{pmatrix}.$$

$$F = \begin{pmatrix} 0 & -1 & -1 \\ 1 & 2 & 1 \\ 1 & 1 & 2 \end{pmatrix},$$

$$\begin{aligned} |F - \lambda I| &= \begin{vmatrix} -\lambda & -1 & -1 \\ 1 & 2 - \lambda & 1 \\ 1 & 1 & 2 - \lambda \end{vmatrix} \\ &= -\lambda(2 - \lambda)^2 - 2 + 2 - \lambda + \lambda + 2 - \lambda \\ &= (2 - \lambda)[- \lambda(2 - \lambda) + 1] \\ &= (2 - \lambda)(-2\lambda + \lambda^2 + 1) \\ &= (2 - \lambda)(1 - \lambda)^2. \end{aligned}$$

One root is  $\lambda_1 = 1$  with algebraic dimension 2, and  $\lambda_2 = 2$  with algebraic dimension 1.

$$Fv = \lambda v \iff \begin{cases} -y - z &= \lambda x \\ x + 2y + z &= \lambda y \\ x + y + 2z &= \lambda z \end{cases}.$$

$\lambda_1 = 1$ :

$$\begin{cases} -y - z &= x \\ x + 2y + z &= y \\ x + y + 2z &= z \end{cases} \iff \begin{cases} -y - z &= x \\ x + y + z &= 0 \end{cases}$$

If  $x = 0$ :  $y = -z$ .

If  $x = t \neq 0$ :  $y + z = -t$ . This works for  $y = t, z = -2t$ .

Therefore,

$$E(\lambda_1) = \left\{ \begin{pmatrix} 0 \\ t \\ -t \end{pmatrix}, \begin{pmatrix} t \\ t \\ -2t \end{pmatrix} \right\}.$$

$\lambda_2 = 2$ :

$$\begin{aligned} \begin{cases} -y - z &= 2x \\ x + 2y + z &= 2y \\ x + y + 2z &= 2z \end{cases} &\iff \begin{cases} -y - z &= 2x \\ x + z &= 0 \\ x + y &= 0 \end{cases} \iff \begin{cases} x = -z \\ x = -y \end{cases} \\ &\iff 2x - y + z = 0 \end{aligned}$$

So

$$E(\lambda_2) = \begin{pmatrix} t \\ -t \\ -t \end{pmatrix}.$$

Another way to represent eigenvalues and eigenvectors is

$$AV = V\Lambda,$$

where  $V = [v_1, \dots, v_n]$  is the matrix formed by putting each eigenvector as a column, and

$$\Lambda = \begin{pmatrix} \lambda_1 & & & \\ & \lambda_2 & & \\ & & \ddots & \\ & & & \lambda_n \end{pmatrix}$$

is the diagonal matrix formed by all eigenvalues.

A matrix  $A$  is **diagonalizable** if there exist  $n$  linearly independent eigenvectors. That is, if the matrix  $V$  is invertible:

$$\Lambda = V^{-1}AV.$$

This leads naturally to the **eigen-decomposition** of the matrix,

$$A = V\Lambda V^{-1}.$$

A real matrix,  $U$ , is **orthogonal** if  $U^T U = U U^T = I$ .

**Proposition 2.3.** *The following statements are equivalent:*

- $U^T$  is orthogonal.
- $U^T = U^{-1}$ .
- $|U| = 1$ .
- $U$ 's eigenvectors are orthonormal (the pairwise dot product is 0 and the norm is 1).

**Example 2.6.** Some examples of orthogonal matrices:

- Identity:  $I$
- Permutation of coordinates:  $\begin{bmatrix} 0 & 1 & 0 \\ 1 & 0 & 0 \\ 0 & 0 & 1 \end{bmatrix}$
- Rotation:  $\begin{pmatrix} \cos \theta & -\sin \theta \\ \sin \theta & \cos \theta \end{pmatrix}$ .
- Reflection:  $\begin{pmatrix} \cos \theta & \sin \theta \\ \sin \theta & -\cos \theta \end{pmatrix}$ .

A matrix  $A$  is said to be **positive semi-definite** when it can be obtained as the product of a matrix by its transpose:

$$\exists X | A = X X^T.$$

Positive semi-definite matrices are always symmetric, because

$$A^T = (X X^T)^T = X X^T = A.$$

A symmetric matrix  $A$  is positive semi-definite if all its eigenvalues are non-negative.

**Proposition 2.4.** *Let  $A$  be a positive semi-definite matrix. Then:*

- $0 \leq \lambda_1 \leq \lambda_2 \leq \dots \leq \lambda_n$  and its eigenvectors are pairwise orthogonal when their eigenvalues are different.
- The eigenvalues are composed of real values.
- The multiplicity of an eigenvalue is the dimension of its eigenspace.

In this case, since eigenvectors are orthogonal, it is possible to store all the eigenvectors in an orthogonal matrix. Therefore, the eigen-decomposition of a positive semi-definite matrix,  $A$ , could be

$$A = U\Lambda U^T,$$

with  $U$  an orthogonal matrix.

As a consequence, the eigen-decomposition of a positive semi-definite matrix is often referred to as its diagonalization.

An alternative definition for positive semi-definite matrix is:

$A$  is positive semi-definite if  $x^T A x \geq 0, \forall x$ .

If it is  $x^T A x > 0, \forall x$ , then it is positive definite.

If it is  $x^T A x \leq 0, \forall x$ , then it is negative semi-definite.

If it is  $x^T A x < 0, \forall x$ , then it is negative definite.

The **rank** of a matrix is the dimension of the vector space generated by its columns (or rows). This corresponds to the maximum number of linearly independent columns of  $A$ . A matrix whose rank is equal to its size is called a **full rank matrix**. Only full rank matrices have an inverse.

**Proposition 2.5.** *The sum of the eigenvalues of a matrix is the sum of the elements of its main diagonal. The product of the eigenvalues is equal to the determinant of the matrix.*

We can now define the **Laplacian matrix** for undirected graphs, as

$$L_{ij} = \begin{cases} -1 & , (v_i, v_j) \in E \\ 0 & , (v_i, v_j) \notin E \\ d_i & , i = j \end{cases}$$

or, equivalently,

$$L = D - A,$$

where  $D$  is the degree matrix of  $G$ , and  $A$  its adjacency matrix.

### 2.2.1 Exercises

1. What could you say about these matrices?

(a)  $A = \begin{pmatrix} -1 & \frac{3}{2} \\ 1 & -1 \end{pmatrix}$ ,  $\det(A) = -\frac{1}{2}$ ,  $A$  is invertible. Its eigenvalues are  $\lambda_1 = -1 + \frac{\sqrt{6}}{2}$  and  $\lambda_2 = -1 - \frac{\sqrt{6}}{2}$ , with  $v_{\lambda_1} = \begin{pmatrix} \frac{\sqrt{6}}{2}t \\ t \end{pmatrix}$  and  $v_{\lambda_2} = \begin{pmatrix} -\frac{\sqrt{6}}{2}t \\ t \end{pmatrix}$ .

(b)  $B = \begin{pmatrix} -1 & \frac{3}{2} \\ \frac{2}{3} & -1 \end{pmatrix}$ . The second row is equal to the first row multiplied by  $-\frac{2}{3}$ . Therefore, it is not invertible.

(c)  $I$ : its determinant is 1. It is symmetric, orthogonal, its own inverse. Triple eigenvalue 1, with eigenspace the whole space.

2. Show that  $A^n = X\Lambda X^{-1}$ .

First, this is only true if  $A$  is diagonalizable. If that is the case, then we can proceed by induction on  $n$ :  
 $n = 1$ : Obvious.

$n = 2$ :

$$A^2 = (X\Lambda X^{-1})^2 = X\Lambda X^{-1}X\Lambda X^{-1} = X\Lambda^2 X^{-1}.$$

Suppose it is true for  $n - 1$ :

$$A^{n-1} = X\Lambda^{n-1}X^{-1}.$$

Then, for  $n$ , we have:

$$A^n = AA^{n-1} = X\Lambda X^{-1}X\Lambda^{n-1}X^{-1} = X\Lambda^n X^{-1}.$$

3. Find the eigenvalues and unit eigenvectors of  $A^T A$  and  $AA^T$  with  $A = \begin{pmatrix} 1 & 1 \\ 1 & 0 \end{pmatrix}$  the Fibonnaci matrix.

First of all, notice that  $A$  is symmetric, so  $A^T A = AA^T = A^2 = \begin{pmatrix} 2 & 1 \\ 1 & 1 \end{pmatrix}$ .

$\begin{vmatrix} 2-\lambda & 1 \\ 1 & 1-\lambda \end{vmatrix} = (2-\lambda)(1-\lambda) - 1 = 2 - 3\lambda + \lambda^2 - 1 = \lambda^2 - 3\lambda + 1$ . The roots of this polynomial are

$$\lambda = \frac{3 \pm \sqrt{9-4}}{2} = \frac{3 \pm \sqrt{5}}{2}.$$

Now,

$$A^2 v = \lambda v \iff \begin{cases} 2x + y = \lambda x \\ x + y = \lambda y \end{cases} \iff \begin{cases} x = (\lambda - 1)y \end{cases}$$

Therefore

$$E(\lambda_1) = \begin{pmatrix} \frac{1+\sqrt{5}}{2}t \\ t \end{pmatrix}$$

with unit eigenvector  $v_1 = \frac{1}{\sqrt{4-\sqrt{5}}} \begin{pmatrix} \frac{1+\sqrt{5}}{2} \\ 1 \end{pmatrix}$ .

And

$$E(\lambda_2) = \begin{pmatrix} \frac{1-\sqrt{5}}{2}t \\ t \end{pmatrix}$$

with unit eigenvector  $v_2 = \frac{1}{\sqrt{4-\sqrt{5}}} \begin{pmatrix} \frac{1-\sqrt{5}}{2} \\ 1 \end{pmatrix}$ .

4. Without multiplying

$$S = \begin{pmatrix} \cos \theta & -\sin \theta \\ \sin \theta & \cos \theta \end{pmatrix} \begin{pmatrix} 2 & 0 \\ 0 & 5 \end{pmatrix} \begin{pmatrix} \cos \theta & \sin \theta \\ -\sin \theta & \cos \theta \end{pmatrix},$$

find the determinant, the eigenvalues and eigenvectors. Why  $S$  is positive definite?

We have  $S = U\Lambda U^T$  with  $U$  orthogonal. Therefore, the eigenvalues of  $S$  are 2 and 5. Its determinant is 10. The eigenvectors are the eigenvectors of  $\Lambda$  rotated as well, that is:

$$V = \begin{pmatrix} \cos \theta & -\sin \theta \\ \sin \theta & \cos \theta \end{pmatrix} \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix} \begin{pmatrix} \cos \theta & \sin \theta \\ -\sin \theta & \cos \theta \end{pmatrix}.$$

$S$  is positive definite because

$$xSx^T = x \begin{pmatrix} \cos \theta & -\sin \theta \\ \sin \theta & \cos \theta \end{pmatrix} \begin{pmatrix} 2 & 0 \\ 0 & 5 \end{pmatrix} \begin{pmatrix} \cos \theta & \sin \theta \\ -\sin \theta & \cos \theta \end{pmatrix} x^T,$$



now note that

$$\left( x \begin{pmatrix} \cos \theta & -\sin \theta \\ \sin \theta & \cos \theta \end{pmatrix} \right)^T = \begin{pmatrix} \cos \theta & \sin \theta \\ -\sin \theta & \cos \theta \end{pmatrix} x^T,$$

so

$$xSx^T = y \begin{pmatrix} 2 & 0 \\ 0 & 5 \end{pmatrix} y^T \geq 0,$$

because  $\begin{pmatrix} 2 & 0 \\ 0 & 5 \end{pmatrix}$  is positive semi-definite (symmetric with positive eigenvalues).

5. For what numbers  $c$  and  $d$  are the following matrices positive definite?

(a)  $A = \begin{pmatrix} c & 1 & 1 \\ 1 & c & 1 \\ 1 & 1 & c \end{pmatrix}$ : all principal minors must be positive. That is:

- $c > 0$ .
- $\begin{vmatrix} c & 1 \\ 1 & c \end{vmatrix} = c^2 - 1 > 0$ . Combined with the previous one, this is  $c > 1$ .
- $\begin{vmatrix} c & 1 & 1 \\ 1 & c & 1 \\ 1 & 1 & c \end{vmatrix} = c^3 + 2 - 3c$ . Roots: 1,  $\frac{1}{2} \begin{vmatrix} 1 & 0 & -3 & 2 \\ & 1 & 1 & -2 \\ & 1 & 1 & -2 & 0 \end{vmatrix}$ , and we have  $c^2 + c - 2$ , with roots  $c = \frac{-1 \pm \sqrt{5}}{2}$ . We are only interested in the interval  $(1, \infty)$ , in which  $c^3 - 3c + 2 > 0$ .

Therefore, it is  $c > 1$ .

(b)  $B = \begin{pmatrix} 1 & 2 & 3 \\ 2 & d & 4 \\ 3 & 4 & 5 \end{pmatrix}$ :

- $1 > 0$ .
- $\begin{vmatrix} 1 & 2 \\ 2 & d \end{vmatrix} = d - 4 > 0 \iff d > 4$ .
- $\begin{vmatrix} 1 & 2 & 3 \\ 2 & d & 4 \\ 3 & 4 & 5 \end{vmatrix} = 5d + 24 + 24 - 9d - 16 - 20 = -4d + 12 > 0 \iff -4d > -12 \iff d < 3$ .

Therefore, there is no value for  $d$  for which  $B$  is positive.

6. Show that if  $\lambda_1, \lambda_2, \dots, \lambda_n$  are the eigenvalues of a matrix  $A$ , then  $A^m$  has as eigenvalues  $\lambda_1^m, \lambda_2^m, \dots, \lambda_n^m$ .

Induction on  $m$ .

$m = 1$ : Obvious.

$m = 2$ : Let  $v_i$  be the eigenvector associated to  $\lambda_i$ , then

$$A^2 v_i = A(Av_i) = A(\lambda_i v_i) = \lambda_i A v_i = \lambda_i^2 v_i,$$

so  $\lambda_i^2$  is an eigenvalue of  $A^2$ , with associated eigenvector  $v_i$ .

Suppose it is true for  $m - 1$ , then, for  $m$ :

$$A^m v_i = A(A^{m-1} v_i) = A(\lambda_i^{m-1} v_i) = \lambda_i^{m-1} A v_i = \lambda_i^m v_i,$$

and we have the result.

7. What is the determinant of any orthogonal matrix?

If  $U$  is orthogonal, then  $UU^T = I$ . Then,

$$1 = |I| = |UU^T| = |U| |U^T| = |U|^2.$$

Therefore,  $|U| = \pm 1$ .

8. For an undirected graph, both the adjacency matrix and the Laplacian matrix are symmetric. Show that the Laplacian matrix is positive semi-definite.

### 3 Random Walks on Graphs

#### 3.1 First Perron-Frobenius Theorem

**Definition 3.1.** A matrix,  $A$ , is **positive** if  $A_{ij} > 0, \forall i, j$ . Similarly, it is **non-negative** if  $A_{ij} \geq 0, \forall i, j$ . Similar definitions apply for negative and non-positive matrices.

*Remark 3.1.* Observe that it is not the same for a matrix to be positive as to be positive semi-definite.

The Perron-Frobenius theorem for non-negative matrices leads to the characterization of non-negative primary eigenvectors. This is useful in stationary distributions, such as those of Markov chains and the famous Google's page rank algorithm.

**Theorem 3.1. Perron-Frobenius Theorem for positive matrices**

If  $A$  is a positive matrix, then:

- $\exists \lambda^* > 0, v^* > 0, \|v^*\|_2 = 1$  such that  $A \cdot v = \lambda^* v^*$  ( $v^*$  is a right column eigenvector).
- $\exists \lambda^* > 0, w > 0, \|w\|_2 = 1$  such that  $w \cdot A = \lambda^* w$  ( $w$  is a left row eigenvector).
- For any other eigenvalue,  $\lambda$ , it holds,  $|\lambda| < \lambda^*$  ( $\lambda^*$  is a dominant eigenvalue, called the **Perron eigenvalue**).
- $\lambda^*$  is unique and  $v^*$  is unique (the only vector of unit length associated to  $\lambda^*$ ).

**Definition 3.2.** A non-negative matrix  $A$  is:

- **Irreducible** if,  $\forall i, j, \exists k \in \mathbb{N}^*$  such that  $A_{i,j}^k > 0$ .
- **Primitive** if,  $\exists k \in \mathbb{N}^*$  such that  $\forall i, j, A_{i,j}^k > 0$ .

**Theorem 3.2. Perron-Frobenius Theorem for non-negative matrices**

If  $A$  is a non-negative matrix, then:

- $\exists \lambda^* > 0, v^* \geq 0, \|v^*\|_2 = 1$  such that  $A \cdot v = \lambda^* v^*$  ( $v^*$  is a right column eigenvector).
- $\exists \lambda^* > 0, w \geq 0, \|w\|_2 = 1$  such that  $w \cdot A = \lambda^* w$  ( $w$  is a left row eigenvector).
- For any other eigenvalue,  $\lambda$ , it holds,  $|\lambda| \leq \lambda^*$  ( $\lambda^*$  is a dominant eigenvalue, called the **Perron eigenvalue**).
- If  $A$  is irreducible, then the vector  $v^*$  is unique and it holds  $v^* > 0$ .
- If  $A$  is primitive, then the eigenvalue  $\lambda^*$  is unique.

Note now that a graph,  $G = (V, E)$ , with adjacency matrix  $A$ , then:  $G$  is connected  $\iff \forall 1 \leq i, j \leq |V|, \exists k \in \mathbb{N}^*$  such that  $A_{i,j}^k > 0$ . This means that the adjacency matrix of connected graphs is irreducible.

Now, if a graph is  $k$ -connected, i.e., there is a  $k$ -path between all nodes, then its adjacency matrix is primitive. One sufficient condition for a graph to be  $k$ -connected is being connected and having  $A_{ii} > 0$  for some  $i$ .

#### 3.2 Random Walks on Graphs

A **random walk** on a graph,  $G = (V, E)$ , is a random process that starts from some vertex  $v_i$ , and repeatedly moves to a neighbour  $v_j$  chosen at random (for example with uniform distribution). The random walk,  $\xi_t$ , is

therefore a random variable describing the position of a random walk after  $t$  steps. The probability of going from node  $i$  to node  $j$  is the **transition probability**,

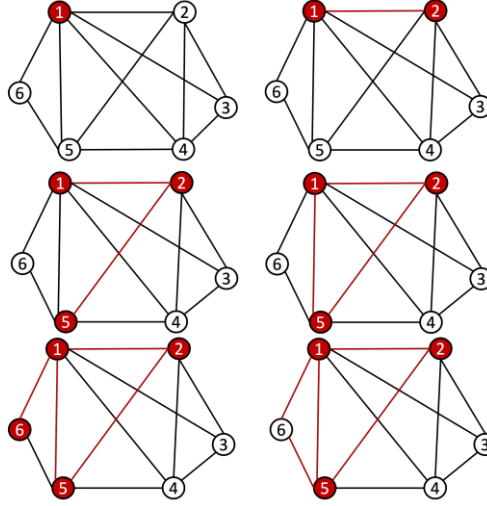
$$P_{ij} = P(\xi_{t+1} = j | \xi_t = i).$$

The sequence of nodes can be regarded as a Markov chain, i.e. a discrete time stochastic process, where the position  $\xi_0$  is the initial state, according to the **init distribution**,  $P^0$ , and from this point the next state only depends on the current state. The  $t$ -step **transition probability** is

$$P_{ij}^t = P(\xi_t = j | \xi_0 = i).$$

Some examples are the path traced by a molecule in a liquid or a gas (Brownian motion), the price of a fluctuating stock, the financial status of a gambler, etc. The term random walk was first introduced by Karl Pearson in 1905.

The following is a basic visual example of a random walk on a graph:



Note that we can express the transition probability  $P_{ij}$  in a matrix  $P$ . This matrix is the **transition probabilities matrix**, and it is **row-stochastic** or **row-Markov**, meaning,

$$P_{ij} \geq 0, \forall i, j, \text{ and } \sum_j P_{i,j} = 1, \forall i.$$

This implies that

$$P \cdot \mathbf{1} = P \cdot \begin{pmatrix} 1 \\ \vdots \\ 1 \end{pmatrix} = \begin{pmatrix} 1 \\ \vdots \\ 1 \end{pmatrix}.$$

This means that  $\begin{pmatrix} 1 \\ \vdots \\ 1 \end{pmatrix}$  is an eigenvector and 1 is an eigenvalue. 1 is the largest eigenvalue because

$$\|Pv\|_1 \leq \|v\|_1,$$

so, for an eigenvalue  $\lambda$ ,

$$|\lambda| \|v\|_1 = \|\lambda v\|_1 = \|Pv\|_1 \leq \|v\|_1,$$

so  $|\lambda| \leq 1$ .

From the Perron-Frobenius theorem for non-negative matrices, we know that:

- $v^* = \begin{pmatrix} 1 \\ \vdots \\ 1 \end{pmatrix}$  is a right Perron eigenvector for  $P$ .
- $|\lambda| \leq \lambda^* = 1$  is a Perron eigenvalue.
- There exists a left Perron eigenvector  $\pi P = \pi$ .
- If  $P$  is irreducible, the vector  $\pi$  is unique.
- If  $P$  is primitive, the eigenvalue 1 is unique (there are no complex eigenvalues with norm 1).

### 3.2.1 The Stationary Distribution

Let  $\pi^t$  be the row vector giving the probability distribution of  $\xi_t$ , that is,  $\pi_i^t$  is the probability that the random walk is at node  $i$  at time  $t$ . Therefore, we can write

$$\pi^{t+1} = \pi^t P,$$

which, applied recursively, leads to

$$\pi^{t+1} = \pi^0 P^{t+1}.$$

Or, we can take limits

$$\lim_t \pi^{t+1} = \lim_t \pi^t P.$$

If this limit exists,  $\lim_t \pi^t = \pi$ , then

$$\pi = \pi P.$$

Convergence is ensured if  $P$  is irreducible.

**Example 3.1.** The following example does not converge:

$$P = \begin{pmatrix} 0 & 1 \\ 1 & 0 \end{pmatrix}$$

A common way to perform random walks on graphs is with the uniform probability. That is,

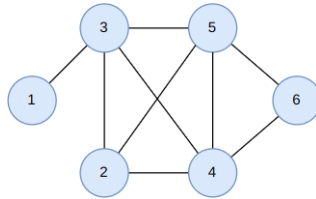
$$P_{ij} = P(\xi_{t+1} = j | \xi_t = i) = \begin{cases} \frac{1}{d_i} & \text{if } (i, j) \in E, \\ 0 & \text{otherwise,} \end{cases}$$

where  $d_i$  is the degree of node  $i$ . Equivalently,

$$P_{ij} = \frac{A_{ij}}{\sum_{j \in V} A_{ij}} = \frac{A_{ij}}{d_i} = D_{ii}^{-1} A_{ij}.$$

The random sequence of vertices  $\xi_0, \xi_1, \dots, \xi_t, \xi_{t+1}, \dots$  visited on  $G$  is a Markov Chain with state space  $V$  and matrix transition probabilité  $P = D^{-1}A$ .

**Example 3.2.** Given the graph:



The transition matrix for the uniform distribution is:

$$P = \begin{pmatrix} 0 & 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & \frac{1}{3} & \frac{1}{3} & \frac{1}{3} & 0 \\ \frac{1}{4} & \frac{1}{4} & 0 & \frac{1}{4} & \frac{1}{4} & 0 \\ 0 & \frac{1}{4} & \frac{1}{4} & 0 & \frac{1}{4} & \frac{1}{4} \\ 0 & \frac{1}{4} & \frac{1}{4} & \frac{1}{4} & 0 & \frac{1}{4} \\ 0 & 0 & 0 & \frac{1}{2} & \frac{1}{2} & 0 \end{pmatrix}$$

### 3.2.2 Balance Condition

A probability distribution  $\pi$  satisfies the **balance condition** if

$$\pi_i P_{ij} = \pi_j P_{ji}, \forall i, j \in V.$$

If  $\pi$  satisfies the balance condition, then it is the stationary distribution for the undirected graph. To see this, notice that the balance condition can be rewritten as

$$\pi_i \frac{A_{ij}}{d_i} = \pi_j \frac{A_{ji}}{d_j}.$$

Since the graph is considered without direction,  $A_{ij} = A_{ji}$ , and then

$$\frac{\pi_i}{d_i} = \frac{\pi_j}{d_j} = c,$$

where  $c$  is a constant, for all  $i, j$ . Now, we know that  $\sum_i \pi_i = 1$ , so

$$1 = \sum_i \pi_i = \sum_i \frac{\pi_j}{d_i} d_i = \sum_i c d_i = c \sum_i d_i.$$

Therefore

$$\sum_i d_i = \frac{1}{c}.$$

Finally, it must be

$$\pi_i = d_i c = \frac{d_i}{\sum_j d_j} = \frac{d_i}{2|E|}.$$

In this case:

$$(\pi P)_i = \sum_j \pi_j P_{ji} = \sum_j \pi_j \frac{1}{d_j} A_{ji} = \sum_j c A_{ji} = c \sum_j A_{ji} = c \sum_j A_{ij} = c d_i = \frac{d_i}{\sum_j d_j} = \pi_i.$$

Therefore, we have seen that the stationary probabilities are proportional to the degrees of the vertices.

In particular, if  $G$  is  $d$ -regular, i.e., all nodes have degree  $d$ , then

$$\pi = \frac{d}{2|E|} = \frac{d}{d \cdot n} = \frac{1}{n}$$

is the uniform distribution. With this setup, a random walk moves along every edge with the same frequency.

The balance condition implies time-reversibility: the reversed walk is also a Markov chain.

### 3.2.3 Hitting Time

**Definition 3.3.** The **expected hitting time**,  $H_{ij}$ , is the expected number of steps before node  $j$  is reached in a random walk starting at node  $i$ :

$$H_{ij} = \begin{cases} 1 + \sum_k P_{ik} H_{kj} & \text{if } i \neq j, \\ 0 & \text{otherwise.} \end{cases}$$

*Remark.* In general,  $H_{ij} \neq H_{ji}$ , so  $H$  is not symmetric.

*Remark 3.2.*  $H$  follows the triangle inequality

$$H_{ij} \leq H_{ik} + H_{kj}.$$

**Definition 3.4.** The **commute time**,  $C_{ij}$ , is the expected number of steps in a random walk starting at node  $i$ , reaching node  $j$  and coming back to  $i$  again:

$$C_{ij} = H_{ij} + H_{ji}.$$

### 3.2.4 Lazy Random Walk

The lazy random walk is a variation of the random walk, in which the walk stays at the current node with probability  $\frac{1}{2}$ , and continue with the walk with the rest of the probability.

In this case, the transition matrix is

$$P_{ij} = \begin{cases} \frac{1}{2} & \text{if } i = j, \\ \frac{1}{2d_i} & \text{if } (i, j) \in E, \\ 0 & \text{otherwise.} \end{cases}$$

If  $Q$  is the transition matrix for the uniform random walk, then

$$\pi^{t+1} = \pi^t P = \frac{1}{2} \pi^t + \frac{1}{2} \pi^t Q.$$

**Proposition 3.1.** *If the lazy random walk converges and the uniform random walk is irreducible, then it converges to the same stationary distribution as the uniform random walk.*

*Proof.* Let  $Q$  be the transition matrix for the uniform random walk, then, the stationary distribution is

$$\pi = \pi Q.$$

For lazy random walk, say the stationary distribution is  $\pi'$ . Then:

$$\pi' = \frac{1}{2} \pi' + \frac{1}{2} \pi' Q \iff \frac{1}{2} \pi' = \frac{1}{2} \pi' Q \iff \pi' = \pi' Q.$$

Therefore, since  $Q$  is irreducible, the uniqueness of  $\pi$  implies  $\pi' = \pi$ . □

## 3.3 PageRank

The web is very heterogeneous bu nature, and certainly huge. We cannot expect the web graph to be connected. Page and Brin proposed a way to overcome this problem, by ensuring the convergence of random walks on the web graph.

The idea is to fix a positive constant,  $p$ , between 0 and 1, called the **damping factor**, and which represents the probability that a user leaves the current page and goes to a random web.

Therefore, the **page rank transition matrix** is

$$P_g = (1 - p) P + pB,$$

where  $B = \frac{1}{n} \begin{pmatrix} 1 & 1 & \dots & 1 \\ 1 & 1 & \dots & 1 \\ \vdots & \vdots & \ddots & \vdots \\ 1 & 1 & \dots & 1 \end{pmatrix}$ .

$p$  is usually chosen small, like 0.15, modelling a situation in which a surfer will, most of the time, follow the outgoing links and move on to one of the neighbours. A smaller percentage of time, the surfer will dump the current page and choose arbitrarily a different page from the web.

**Proposition 3.2.**  $P_g$  is stochastic.

*Proof.* We need to proof that, for all  $i$ , it holds  $\sum_j P_{g_{i,j}} = 1$ .

$$\begin{aligned} \sum_j P_{g_{i,j}} &= \sum_j (1-p) P_{ij} + p B_{ij} \\ &= (1-p) \sum_j P_{ij} + p \sum_j B_{ij} \\ &= (1-p) \cdot 1 + p \sum_j \frac{1}{n} \\ &= 1-p + p \cdot n \frac{1}{n} \\ &= 1-p + p \\ &= 1. \end{aligned}$$

□

## 4 Centrality Measures

Centrality Measures try to answer the question 'What characterizes an important vertex?'. They define a real-valued function on the vertices of the graph,  $m : V \rightarrow \mathbb{R}$ , that serves to rank the vertices. However, there are many different ways to define such function, leading to different definitions of centrality, such as cohesiveness, ability to transfer information across the network, to influence other nodes, to control information flow, etc.

There are many centrality measures that count the number of paths through a given vertex. These differ in how relevant walks are defined and counted. For example, if we only consider paths of length one, we would be computing degree centrality, while if we allow paths of arbitrary length, we would be computing eigenvalue centrality.

### 4.1 Degree Centrality

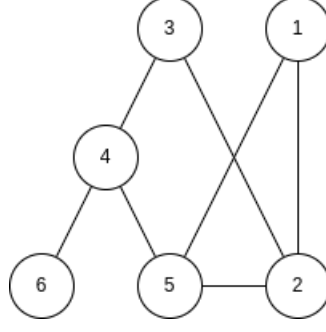
The more neighbours a vertex has, the higher its communication ability is, increasing its importance.

**Definition 4.1.** Given the graph  $G = (V, E)$ , with adjacency matrix  $A$ , the **degree centrality** is computed as

$$D = Au,$$

where  $u = \mathbf{1} \in \mathbb{R}^n$ .

**Example 4.1.** Consider the following graph:



The degree centrality is

$$D = Au = \begin{pmatrix} 0 & 1 & 0 & 0 & 1 & 0 \\ 1 & 0 & 1 & 0 & 1 & 0 \\ 0 & 1 & 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 & 1 & 1 \\ 1 & 1 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 & 0 \end{pmatrix} \begin{pmatrix} 1 \\ 1 \\ 1 \\ 1 \\ 1 \\ 1 \end{pmatrix} = \begin{pmatrix} 2 \\ 3 \\ 2 \\ 3 \\ 3 \\ 1 \end{pmatrix},$$

so the nodes with highest value are nodes (2, 4, 5).

One drawback of this measure, is that it is very likely that several nodes present the same exact value, diffculting an unique ranking of vertices.

## 4.2 Neighbourhood centrality

This measure correspond to the average degree of each vertex neighbours. We could understand this measure as measuring how much a vertex is related to influential vertices.

**Definition 4.2.** Given the graph  $G = (V, E)$ , with adjacency matrix  $A$ , the **neighbourhood centrality** is computed as

$$N = \mathcal{D}^{-1}AD,$$

where  $\mathcal{D}$  is the diagonal matrix where  $\mathcal{D}_{ii} = d_i$  is the degree of vertex  $i$  and  $D$  is the degree centrality. Each vertex' measure is

$$N_v = \frac{\sum_{u \in \mathcal{N}_v} d_u}{d_v}.$$

**Example 4.2.** The neighbourhood centrality of the previous example graph is

$$N = \mathcal{D}^{-1}Au = \begin{pmatrix} \frac{1}{2} & & & & & \\ & \frac{1}{3} & & & & \\ & & \frac{1}{2} & & & \\ & & & \frac{1}{3} & & \\ & & & & \frac{1}{3} & \\ & & & & & 1 \end{pmatrix} \begin{pmatrix} 0 & 1 & 0 & 0 & 1 & 0 \\ 1 & 0 & 1 & 0 & 1 & 0 \\ 0 & 1 & 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 & 1 & 1 \\ 1 & 1 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 & 0 \end{pmatrix} \begin{pmatrix} 2 \\ 3 \\ 2 \\ 3 \\ 3 \\ 1 \end{pmatrix} = \begin{pmatrix} \frac{3+3}{2} \\ \frac{2+2+3}{3} \\ \frac{3+3}{2} \\ \frac{2+3+1}{3} \\ \frac{2+3+3}{3} \\ \frac{3}{1} \end{pmatrix} = \begin{pmatrix} 3 \\ 2.33 \\ 3 \\ 2 \\ 2.66 \\ 3 \end{pmatrix},$$

so the nodes with highest value are nodes (1, 3, 6).

## 4.3 Eigenvector Centrality

A natural extension of degree centrality is to consider all reachable nodes, not just neighbours. Eigenvector centrality measures a node's importance while considering the importance of its neighbours. A high eigenvector centrality means that a node is connected to many nodes that have high scores themselves.



**Definition 4.3.** Given the graph  $G = (V, E)$ , with adjacency matrix  $A$ , the **eigenvector centrality** of node  $v$  is

$$E_v = \frac{1}{\lambda} \sum_{u \in \mathcal{N}_v} A_{vu} E_u,$$

where  $\lambda$  is a parameter. Note that this can be written as

$$E = \frac{1}{\lambda} A E,$$

or

$$A E = \lambda E.$$

This means that  $E$  is an eigenvector of  $A$ , for the eigenvalue  $\lambda$ .

Bonacich suggested that the eigenvector of the largest eigenvalue of  $A$  could make a good network centrality measure.

The eigenvector  $E$  must be non-negative and according to the Perron-Frobenius theorem, the largest  $\lambda$  enforces this property, making it a suitable value.

**Example 4.3.** Let's compute  $E$  for the previous example graph. The matrix  $A$  has as largest eigenvalue  $\lambda = 2.54$ , and the corresponding eigenvector is

$$E = \begin{pmatrix} 2.5 \\ 3.1 \\ 2.2 \\ 2.5 \\ 3.2 \\ 1 \end{pmatrix}.$$

Note that it is usually unfeasible to compute the eigenvalues and eigenvectors. It is more usual to get the vector iteratively as

$$E_k = A \frac{E_{k-1}}{\|E_{k-1}\|}.$$

#### 4.4 PageRank Centrality

Google's PageRank is a variant of the eigenvector centrality, which uses in-degree to award one centrality point for every link a node receives. As we saw, the algorithm is based on a web surfer who is randomly clicking on links, with a certain probability to go to a different place of the web (the damping factor).

Therefore, we define the matrix

$$P_g = (1 - p) P + p B,$$

where  $P_{ij} = \begin{cases} \frac{1}{d_i} & \text{if } j \in \mathcal{N}_i \\ 0 & \text{otherwise} \end{cases}$ , and  $B_{ij} = \frac{1}{n}$ .

Now, we apply the eigenvector centrality to this modified matrix, as

$$P_g E_g = \lambda E_g = E_g,$$

with  $\lambda = 1$  because  $P_g$  is stochastic.

Or, iteratively as

$$E_{g_k} = P_g E_{g_{k-1}}.$$

Note that in this case it is not necessary to normalize the vector at each step, because  $P_g$  is stochastic. A good

$$E_{g_0} \text{ is } E_{g_0} = \begin{pmatrix} \frac{1}{n} \\ \vdots \\ \frac{1}{n} \end{pmatrix}.$$

### 4.5 Katz (or alpha) centrality

The main problem with eigenvector centrality is that it only works well when the graph is strongly connected (so Perron-Frobenius is applicable in its stronger form). Real networks do not usually have this property, specially if they are directed. The vertices that are not in strongly connected components will have value 0.

A way to work around this problem was proposed by Leo Katz. The idea is to give each node a minimum, positive amount of centrality, that it can transfer to other nodes, so:

$$K_v = \alpha \sum_u A_{vu} K_u + \beta,$$

where  $K_v$  is the Katz centrality of node  $v$ ,  $\beta$  is a vector whose elements are all equal to a given positive constant and  $\alpha \in (0, 1)$  is a parameter. Equivalently, this is

$$K = \alpha AK + \beta,$$

so

$$(I - \alpha A) K = \beta,$$

and

$$K = (I - \alpha A)^{-1} \beta.$$

For this to work,  $I - \alpha A$  must be invertible, which happens if and only if  $|I - \alpha A| \neq 0 \iff |\frac{1}{\alpha} I - A| \neq 0$ , so  $\frac{1}{\alpha}$  must not be an eigenvalue of  $A$ . This is ensured if we take  $\frac{1}{\alpha} > \lambda_{max}$ , or  $0 < \alpha < \frac{1}{\lambda_{max}}$ .

An iterative way to compute  $K$  is

$$K = \left( \sum_{k=1}^{\infty} \alpha^k A^k \right) u.$$

The strength of  $\alpha$  decreases at each iteration, acting as attenuation factor.

### 4.6 Clustering Coefficient Centrality

**Triadic closure** is the property among three nodes A, B, and C (representing people, for instance), that if the connections A-B and A-C exist, there is a tendency for the new connection B-C to be formed.

The clustering coefficient measures the proportion of neighbours of each node, that connected to each other.

**Definition 4.4.** Given a graph  $G = (V, E)$ , with adjacency matrix  $A$ , the **clustering coefficient** of node  $v$  is

$$CC_v = \frac{|\{ \{u, v, w\} : (u, v), (v, w), (u, w) \in E \}|}{\binom{d_v}{2}},$$

where the numerator is the number of triangles involving  $v$  and its neighbours, and the denominator is the total number of possible links between  $v$ 's neighbours.

The more densely connected the neighbourhood of  $v$  is, the higher is its clustering coefficient.

**Example 4.4.** The clustering coefficient of the graph example that we've been working with is

$$CC = \begin{pmatrix} \frac{1}{3} \\ \frac{1}{3} \\ \frac{1}{3} \\ 0 \\ 0 \\ \frac{1}{3} \\ \frac{1}{3} \\ \frac{1}{3} \\ 0 \end{pmatrix} = \begin{pmatrix} \frac{1}{3} \\ \frac{1}{3} \\ 0 \\ 0 \\ 0 \\ \frac{1}{3} \\ \frac{1}{3} \\ \frac{1}{3} \\ 0 \end{pmatrix}.$$

## 4.7 Closeness Centrality

Closeness centrality is a measure of how close a node is, on average, to the rest of the nodes, in terms of shortest paths. It measures the average distance between a node  $v$  and all other nodes in the network. Thus, the more central a node is, the closer it is to all other nodes.

**Definition 4.5.** Given a graph  $G = (V, E)$ , the **closeness centrality** of node  $v$  is

$$CL_v = \frac{1}{\sum_{r \neq v} \text{dist}(v, r)}.$$

It can be normalized by the factor

$$CL_v = \frac{N - 1}{\sum_{r \neq v} \text{dist}(v, r)}.$$

An alternative is the **harmonic centrality**, obtained as

$$H_v = \sum_{r \neq v} \frac{1}{\text{dist}(v, r)},$$

with  $\text{dist}(v, r) = 0$  if there is no path from  $v$  to  $r$ .

## 4.8 Betweenness Centrality

A family of betweenness measures are defined to capture a node's importance as a conduct of information flow in the network. This has wide applications in network theory, because in a telecommunications network, a node with higher betweenness centrality would have more control over the network, since more information will pass through that node.

The most well-known betweenness metric measures the number of times a node is on a shortest path between two nodes.

**Definition 4.6.** Given a graph  $G = (V, E)$ , the **betweenness centrality** of node  $v$  is

$$B_v = \sum_{s \neq v \neq t} \frac{\sigma_{s,t}(v)}{\sigma_{s,t}},$$

where  $\sigma_{s,t}$  is the number of shortest path from source node  $s$  to target node  $t$ , and  $\sigma_{s,t}(v)$  is the number of shortest path between these two nodes going through  $v$ .

This measure can be normalized by the number of ordered pairs not including  $v$ :

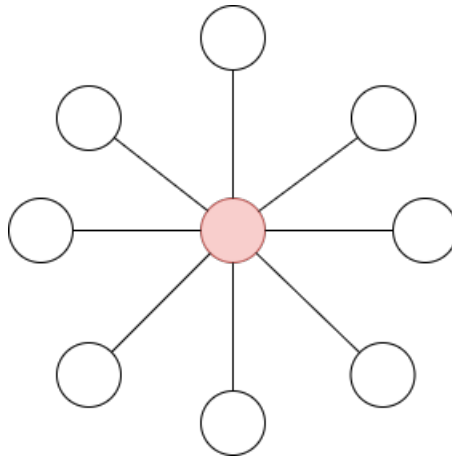
- For directed graphs

$$B_v = \frac{1}{(n-1)(n-2)} \sum_{s \neq v \neq t} \frac{\sigma_{s,t}(v)}{\sigma_{s,t}}.$$

- For undirected graphs

$$B_v = \frac{2}{(n-1)(n-2)} \sum_{s \neq v \neq t} \frac{\sigma_{s,t}(v)}{\sigma_{s,t}}.$$

**Example 4.5.** For the undirected star graph:

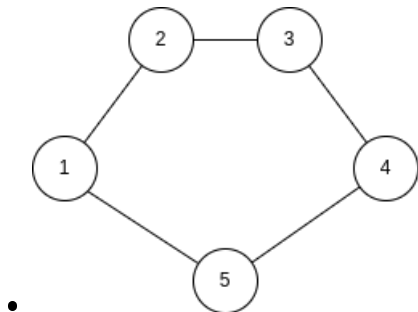


The center vertex has a betweenness of  $\frac{(n-1)(n-2)}{2}$  (or 1, if we normalize it), while the leaves have a betweenness of 0.

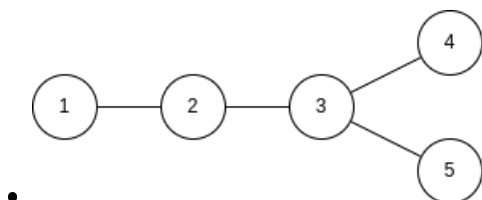
**Exercise 4.1.** What about the following graphs?



$$B = \begin{pmatrix} 0 \\ 3 \\ 6 \\ 4 \\ 3 \\ 0 \\ 6 \end{pmatrix} = \begin{pmatrix} 0 \\ 0.5 \\ 0.67 \\ 0.5 \\ 0 \\ 0 \end{pmatrix}.$$



$$B = \begin{pmatrix} \frac{1}{6} \\ \frac{1}{6} \\ \frac{1}{6} \\ \frac{1}{6} \\ \frac{1}{6} \end{pmatrix}.$$



$$B = \begin{pmatrix} 0 \\ 5 \\ 6 \\ 6 \\ 6 \end{pmatrix} = \begin{pmatrix} 0 \\ 0.5 \\ 0.83 \\ 0 \\ 0 \end{pmatrix}.$$

## 5 MapReduce Computation Model

The advent of big data and the increasing analysis needs favored the design of parallel algorithms, specially in the realm of big data processing pipelines, with a tradeoff between communication costs and degree of parallelism.

MapReduce is a processing paradigm that works on top of distributed environments. More precisely, it was built on top of Google File System (GFS) and Hadoop Distributed File System (HDFS), used to manage large-scale data and to be tolerant to hardware and networks faults. To do this, HDFS splits files into large blocks and distributes them across nodes in a cluster, and MapReduce is the programming model used to manage many large-scale parallel computations.

Basically, the idea is that the data is first splitted, then some operation is done to it, and then it's merged to produce the final results. For this, we will just need to define the **Map** and **Reduce** functions, while the system manages the parallel execution on distributed data and the coordination between them, leading with the possibility that one of the tasks may fail.

### Example 5.1. Word Counter

Consider a text file splitted into partitions  $A, B, C, D$ , across different nodes. We want to count how many times each word appears in the whole document. For this, we can use MapReduce as follows:

1. Map: for each word,  $w$ , in each partition, generate the pair  $(w, 1)$ .
2. Shuffle/sort: collects and groups the pairs by key (word), in order to guarantee that the same key will be processed by the same reduce task. **Shuffling** is the process of redistributing data from Map nodes to Reduce nodes.

In our example, we would have, for each word  $w$ , the pairs  $(w, [1, \dots, 1])$ , with as many 1s as  $w$  appearances.

3. Reduce: for each input  $(w, [1, \dots, 1])$ , output  $(w, N_w)$ , where  $N_w$  is the amount of 1s.

The Map task will typically process many words in one or more chunks. If a word,  $w$ , appears  $m$  times among all chunks assigned to that process, there will be  $m$  key-value pairs  $(w, 1)$  among its output.

To perform the grouping and distribution to the Reduce task, the master controller merges the pairs by key and produces a sequence of  $(w, [1, \dots, 1])$ . Since it knows how many reduce tasks there will be,  $r$ , it will produce  $r$  lists, putting a list in one of  $r$  local files destined to one of the Reduce tasks. Each key is assigned as input to one, and only one, Reduce task.

The Reduce task executes one or more reducers, one per key. The outputs from all reducers are merged into a single final file.

### 5.1 The Map Function

In general, a map function can be defined as a function,  $m_f : \mathbb{E}_1^n \rightarrow \mathbb{E}_2^n$ , where  $\mathbb{E}_i$  is the domain of the input (1) or output (2) and  $f : \mathbb{R} \rightarrow \mathbb{R}$ , that applies  $f$  to each coordinate. That is:

$$m_f([e_1, \dots, e_n]) = [f(e_1), \dots, f(e_n)].$$

For example:

$$m_{\cdot 2}([2, 3, 6]) = [4, 6, 12].$$

In the MapReduce scheme, map is more restrictive, as the function  $f$  must produce a key-value pair. That is, for all  $i = 1, \dots, n$ , it is

$$f(e_i) = (k_i, v_i).$$

For example, in the word counter example:

$$m_f(["a", "b", "a"]) = [("a", 1), ("b", 1), ("a", 1)].$$

## 5.2 Shuffling/Grouping Function

The shuffle function consists in grouping the outputs of the map function by key, so

$$s([(k_1, v_1), \dots, (k_n, v_n)]) = [(k_1, (v_j : k_j = k_1, \forall j = 1, \dots, n)), \dots].$$

Following the previous example:

$$s([("a", 1), ("b", 1), ("a", 1)]) = [("a", [1, 1]), ("b", [1])].$$

## 5.3 Reduce Function

Generally, a reduce function applies to a vector/row, and outputs a single value, applying the aggregation function  $f$ :

$$r_f([v_1, \dots, v_n]) = f(v_1, \dots, v_n).$$

In MapReduce, reduce applies to each output of the shuffle function with the same key:

$$r_f((k, [v_1, \dots, v_n])) = [(k'_1, f(v_1, \dots, v_n)), \dots, (k'_m, f(v_1, \dots, v_n))].$$

Following the previous example:

$$r_{sum}([("a", [1, 1]), ("b", [1])]) = [("a", 2), ("b", 1)].$$

A **MapReduce pipeline** can be a composition of different  $r_{f_r} \circ s \circ m_{f_m}$ .

The process is illustrated below:



## 5.4 MapReduce Execution Model

Whenever we launch the execution of a MapReduce pipeline, the following happens:

- The user program forks a master controller process and some number of worker processes at different computer nodes.
- The amster creates some number of map tasks and some number of reduce tasks. It assigns the tasks to worker processes by taking into account the co-location.
- A worker handles either map tasks (a **map worker**) or reduce tasks (a **reduce worker**), but not both.
- A worker process reports to the amster when it finishes a task, and a new task is scheduled by the master for that worker process.

- The master keeps track of the status of each map and reduce task (idle, executing, or completed).



#### 5.4.1 Coping with Node Failures

If the master node fails, the entire MapReduce job must be restarted.

If a worker node fails, it would be detected and managed by the master, since it periodically pings the worker processes. All the map tasks assigned to this worker have to be redone in this case.

#### 5.4.2 Algorithms by MapReduce

This paradigm is not a solution to every problem, and in fact it only makes sense when files are very large, and rarely outdated. Its original purpose was to execute very large matrix-vector multiplications.

### 5.5 Use-Case: Matrix-Vector Multiplication by MapReduce

Let  $M$  be a  $n \times n$  squared matrix and  $V$  a vector of size  $n$ . Their product,

$$W = MV,$$

is defined by

$$w_i = \sum_{j=1}^n m_{ij}v_j.$$

We can store  $M$  and  $V$  in a file in *HDFS* as triples  $((i, j), m_{ij})$  for  $M$  and pairs  $(j, v_j)$  for  $V$ <sup>1</sup>. Now we can compute the computation by MapReduce as:

- Map: for each  $((i, j), m_{ij})$  and  $(j, v_j)$ , it returns  $(i, m_{ij}v_j)$ .
- Reduce: simply sums all the values for each key  $i$ , producing the pair  $(i, w_i)$ .

For this to work, all the pairs from  $V$  must be available in all chunks ( $V$  cannot be stored distributely).

More concretely, we can define the functions:

---

<sup>1</sup>This way is very efficient for sparse matrices.

```

1 map(key, val):
2   i = key(1)
3   j = key(2)
4   for (j2, v) in V:
5     if j == j2:
6       emit(i, val*v)
7
8 reduce(key, val):
9   sum = 0
10  for v in val
11    sum += v
12  emit(key, sum)

```



Now, if  $n$  is large,  $V$  might not fit in main memory of a worker node, and a large number of disk accesses may be required. We can improve the approach by distributing  $V$  and refining the algorithm as follows:

- We divide the matrix into vertical stripes of equal width, and the vector in strips of the same size:



Here, the size of  $M_k$  is  $n \times n_k$  and the size of  $V_k$  is  $n_k$ , so that the product  $M_k \cdot V_k$  can be performed, outputting a vector of size  $n$ .

- Each map task is assigned a chunk from one of the matrix stripes and gets the entire corresponding stripe of the vector.
- The final result would be

$$W = MV = \sum_{k=1}^K M_k \cdot V_k,$$

where we apply the previously explained algorithm to each sub-multiplication step.

### 5.5.1 Matrix Multiplication

This approach can be extended to matrix multiplication. Now, let  $M$  be a matrix of size  $n_1 \times n_2$  and  $N$  a matrix of size  $n_2 \times n_3$ , the product  $P = MN$  is a matrix of size  $n_1 \times n_3$ , where

$$p_{ik} = \sum_{j=1}^{n_2} M_{ij} N_{jk}.$$

The matrices are stored as  $(M, (i, j), m_{ij})$  and  $(N, (j, k), n_{jk})$ .



- Map 1: transform  $(M, (i, j), m_{ij})$  into  $(j, (M, i, m_{ij}))$  and  $(N, (j, k), n_{jk})$  into  $(j, (N, k, n_{jk}))$ .
- Reduce 1: for each key,  $j$ , produces the key-value pair  $((i, k), m_{ij}n_{jk})$ .
- Map 2: the identity.
- Reduce 2: for each key,  $(i, k)$ , produce the sum of the list of values associated to this key,  $((i, k), \sum_j m_{ij}n_{jk})$ .

In addition,  $M$  could be divided into  $K$  vertical stripes of size  $(n_1, n_k)$  and  $N$  into  $K$  horizontal stripes of size  $(n_k, n_3)$ , where  $\sum_k n_k = n_2$ . In this setup, we can apply the algorithm to compute each  $M_k \cdot N_k$  and then sum them all.

The functions can be defined more precisely as:

```

1 map_1(T, (i, j), T_ij):
2     emit(j, (T, i, T_ij))
3
4 reduce_1(key, val):
5     for v in val:
6         for w in val:
7             if v(1) == M and w(1) == N:
8                 i = v(2)
9                 M_ij = v(3)
10                k = w(2)
11                N_jk = w(3)
12                emit((i, k), M_ij*N_jk)
13
14 map_2(key, val):
15     emit(key, val)
16
17 reduce_2(key, val):
18     sum = 0
19     for v in val:
20         sum += v
21     emit(key, sum)

```

## 5.6 Relational Algebra by MapReduce

### 5.6.1 Selection

Let  $R(A_1, \dots, A_n)$  be a relation stored as a file in HDFS. The elements of this file are the tuples of  $R$ . The selection operator,  $\sigma_C(R)$  can be defined using MapReduce as:

- Map: for each tuple in  $R$ ,  $t$ , test if  $t$  satisfies  $C$ . If it does, produce the key-value pair  $(t, t)$ .
- Reduce: the identity.

### 5.6.2 Projection

For the projection,  $\pi_A(R)$ , we can do:

- Map: for each tuple in  $R$ ,  $t$ , construct a tuple  $t'$  by removing the attributes that are not in  $A$ . Output  $(t', t')$ .
- Reduce: for each key,  $t'$ , produced by the map tasks, there will be one or more key-value pairs  $(t', t')$ . The reduce function turns  $(t', [t', t', \dots, t'])$  into  $(t', t')$  so it produces exactly one pair.

### 5.6.3 Join

$R(A) \bowtie_B S(C)$  with  $A, B, C$  sets of attributes satisfying  $B \subset A, B \subset C$ , can be implemented with MapReduce as:

- Map: for each tuple  $(a, b) \in R$ , produce the key-value pair  $(b, (R, a))$ . For each tuple  $(c, b) \in S$ , produce the key-value pair  $(b, (S, c))$ .
- Reduce: for each key,  $b$ , output as many pairs as needed,  $(b, [(R, a), (S, c)])$ .

### 5.6.4 Aggregation

The aggregation operator,  $\gamma_{A, \theta(B)}(R)$ , where  $A \cup B$  is the set of attributes of  $R$ , and  $A \cap B = \emptyset$ , can be defined with MapReduce as:

- Map: for each tuple,  $t$ , produce  $(a, b)$ , where  $a$  is the  $A$  part of  $t$ , and  $b$  is the  $B$  part.
- Reduce: each key represents a group, so we apply  $\theta$  to the list  $[b_1, \dots, b_n]$  associated to each value  $a$ . We output  $(a, x)$ , where  $x = \theta(b_1, \dots, b_n)$ .

## 5.7 Some Issues of MapReduce

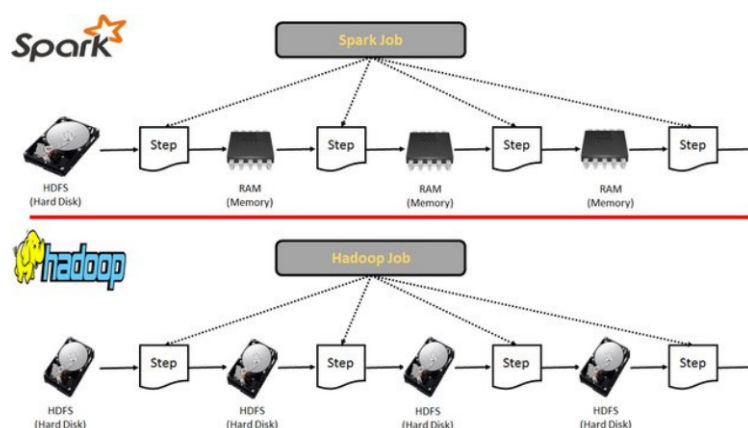
- **Locality:** input data is stored on local disks of machines in the cluster. Each file is divided into blocks of 64MB, each of which is stored several times, as replicas, on different machines. MapReduce master node takes the location information of the input files into account, and attempts to schedule a map task on a machine that contains the needed replica. If this fails, it tries to schedule a map task in a machine that is near to one that has a replica.
- **Granularity:** the map and reduce steps are divided into  $M$  and  $R$  pieces.  $M$  and  $R$  should be much larger than the number of workers. Each worker can perform different tasks, improving dynamic load balancing and speeding up recovery when a worker fails. Some practical bounds on how large these values should be say that the master should take  $M + R$  scheduling decisions and keep  $M \times R$  states in memory.
- **Refinements:** partitioning input data using different functions according to the problem to be solved.
- **Ordering guarantees:** the intermediate key-value pairs are generally processed in increasing key order, to make it easy to generate a sorted output file per partition. However, this is not guaranteed.

## 6 Spark Parallel Computing Framework

Nowadays, data is growing faster than processing speeds, and so the only possible solution is to parallelize on large clusters.

**Apache Spark** is an open source implementation of a framework for large-scale data processing, providing an interface for programming clusters with implicit data parallelism and fault tolerance. It extends a programming language with read-only data structure distributed over a cluster of machines, the **Resilient Distributed Datasets (RDDs)**, maintained in a fault-tolerant way. RDDs were developed in 2012 in response to limitations in Hadoop's MapReduce, which forces a particular linear dataflow as a sequence of HDFS reads and writes.

Spark is up to 100 times faster than traditional Hadoop thanks to its in-memory data processing:



## 6.1 Spark's Engine Properties and Components

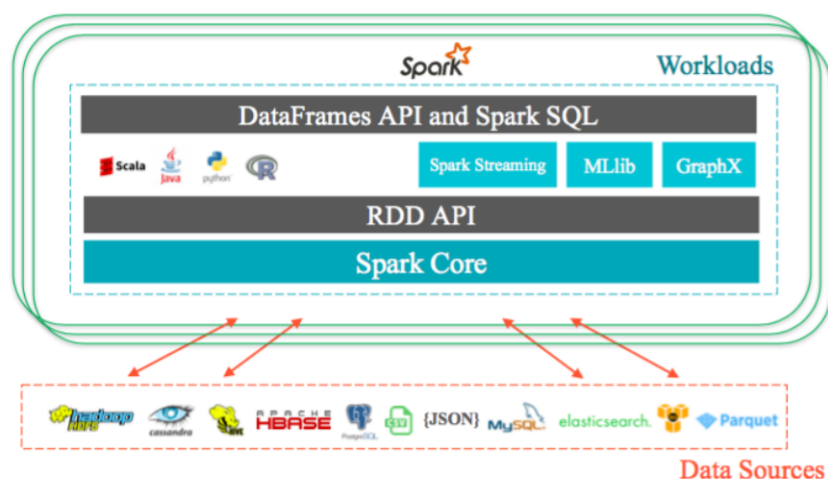
Spark was originally written in Scala, a high level language for JVM. There are APIs for Java, Scala, Python, R,...

The Dataframe API was released as an abstraction on top of the RDD, as well as packages like MLlib or GraphX, that can be used for machine learning and graph analytics. These APIs facilitate the implementation of both iterative algorithms and interactive or exploratory data analysis.

Spark requires a cluster manager and a distributed storage system:

- Cluster management: Spark supports standalone, native Spark clusters, where we can launch a cluster either manually or using scripts, or we can use Hadoop YARN, Apache Mesos, or Kubernetes.
- Distributed storage: Spark can interface with a wide variety of distributed databases, like Alluxio, HDFS, MapR-FS, Cassandra,...

Spark also supports different sources of data, in different formats and from different databases. All these relationships are shown below:



## 6.2 Spark's Execution Architecture

Data is splitted into partitions or blocks, and the driver assigns tasks to each worker, which reads a HDFS block, and has a cache. Each worker process and cache data, if necessary, sending the results to the driver when done.



### 6.3 Resilient Distributed Datasets (RDDs)

RDDs are immutable and distributed collections of objects, spread across a cluster, and stored in RAM or disk (when they are persistent). They are statically typed, i.e., `RDD[T]` has objects of type `T`. The types can be of any type of Python, Java, or Scala objects, including user-defined classes.

RDDs are built via parallel **transformations** and computed via parallel **actions** on distributed datasets, executed lazily. For instance, RDDs are splitted into multiple partitions, which may be computed on different nodes of a cluster.

- **Transformation:** operation on an RDD that returns a new RDD. They are computed lazily, only after an action is called.
- **Action:** operation on an RDD that returns a final result which is not another RDD. When an action is called, all the transformations prior to it are executed in the same order they were defined. Each time a new action is called, the entire RDD must be computed from scratch, but the user can decide to persist intermediate result (caching).

Inside Apache Spark, the workflow is managed as a directed acyclic graph (DAG). Nodes represent RDDs while edges represent the operations executed on the RDDs. Spark keeps track of the set of dependencies between different RDDs. This is called the **lineage graph**.

**Example 6.1.** Python Example: first line mentioning Python

```

1 # spark context creation
2 sc = pyspark.SparkContext(...)
3
4 # creating an RDD of strings with textFile()
5 lines = sc.textFile("README.txt")
6
7 # Transformations:
8 # Construct a new RDD from a previous one, one common transformation is filtering data that
9 # matches a predicate
10 pythonLines = lines.filter(lambda line: "Python" in line)
11
12 # Actions:
13 # Compute a result based on an RDD, and either return it to the driver program or save it to
14 # an external storage system
15 pythonLines.first()

```

Some notes:

- Once a `SparkContext`, `sc`, is set, it is used to build RDDs. The driver program manages a number of workers, and different workers on different machines might count lines in different ranges of the file.
- `filter()` does not mutate the existing input RDD. Instead, it returns a pointer to a new RDD.
- Many transformations work on one element at a time, but this is not true for all transformation, like `union`.

### 6.3.1 RDDs Operations

- `map()` is used from different purposes, from fetching a large collection to just squaring numbers. It transforms an RDD of length  $N$  into another RDD of length  $N$ , applying a function to each element in the RDD. For example:

```

1 nums = sc.parallelize([1,2,3,4])
2 squared = nums.map(lambda x: x*x).collect()
3
4 for num in squared:
5     print(num)
6
7 # 1 4 9 16

```

- `flatMap()` transforms an RDD of length  $N$  into a collection of  $N$  collection, then flattens these into a single RDD of results, applying a function to each collection. For example:

```

1 def tokenize(sentence):
2     return [word for word in sentence]
3
4 rdd = sc.parallelize(["coffee panda"])
5
6 # Using map
7 map_result = rdd1.map(tokenize).collect()
8
9 for token in map_result:
10     print(token)
11
12 # Output: ["coffee", "panda"]
13
14 # Using flatMap
15 flatmap_result = rdd1.flatMap(tokenize).collect()
16
17 for token in flatmap_result:
18     print(token)
19
20 # Output: "coffee", "panda"

```

- `reduce()` is the most common action on basic RDDs. It operates on two elements of the type of the RDD and returns a new element of the same type. For example:

```

1 data = sc.parallelize([1,2,3])
2
3 data.reduce(lambda x, y: x+y) # Result is 6

```

- `reduceByKey()` operates on RDD of key-value pairs. It runs several parallel reduce operations, for each key, where each operations combines values that have the same key. It returns a new RDD consisting of each key and the reduced value for that key. For example:

```

1 pets = sc.parallelize([("cat", 3), ("dog", 2), ("cat", 1)])
2
3 pets.reduceByKey(lambda x, y: x+y)
4
5 # Result: {("cat", 4), ("dog", 2)}

```

- `groupByKey()` and `sortByKey()` return the RDD with the values grouped or sorted by the keys. For example:

```

1 pets = sc.parallelize([("cat", 3), ("dog", 2), ("cat", 1)])
2
3 pets.groupByKey() # Result: {("cat", [3, 1]), ("dog", [2])}
4
5 pets.sortByKey() # Result: {("cat", 3), ("cat", 1), ("dog", 2)}

```

- `distinct()` produces a new RDD with only distinct items.

- *union(otherRDD)* produces a new RDD consisting of the data from both sources, removing duplicates.
- *RDD.subtract(otherRDD)* produces a new RDD consisting in all values present in *RDD* but not in *otherRDD*.
- *RDD.cartesian(otherRDD)* returns of possible pairs  $(a, b)$  where  $a \in RDD$  and  $b \in otherRDD$ . Note that this operation is very costly.
- *RDD.innerJoin(otherRDD)* returns only keys that are present in both pairs to the output RDD.
- *RDD.leftOuterJoin(otherRDD)* and *RDD.rightOuterJoin(otherRDD)* join the two RDDs together by key, allowing one of them to miss the key (left or right). For example:

```

1 rdd = sc.parallelize([(1,2),(3,4),(3,6)])
2 other = sc.parallelize([3, 9])
3
4 rdd.join(other) # result: {(3, (4,9)), (3, (6,9))}
5
6 rdd.leftOuterJoin(other) # result: {(1, (2, None)), (3, (4,9)), (3, (6,9))}

```

### 6.3.2 RDDs Actions

- *collect()* is used to retrieve the entire RDD. It is useful if it filters RDDs down to a very small size to deal with it locally at the driver. The retrieved dataset must fit in memory in a single machine.
- *take(N)* is used to retrieve a small number of elements in the RDD at the driver program and then iterate them over them locally.
- *top(N)* is used to extract the top elements.
- *takeSample(withReplacement, N, seed)* allows to take a sample of the data either with out without replacement.
- *foreach()* performs computations on each element in the RDD without bringing it back locally.
- *count()* returns a count of the elements.
- *countByValue()* returns a map of each unique value and its count.

### 6.3.3 Caching

Spark RDDs are lazily evaluated, and sometimes we use the same RDD multiple times. Naively, Spark will recompute the RDD and all of its dependencies each time we call an action on the RDD. To avoid this, we can ask Spark to persist data using *persist()*.

Notice that calling *persist()* does not force the evaluation of the RDD.

If we cache too much data, Spark will automatically delete old partitions. For the memory-only storage levels, it will recompute these partitions the next time they are accessed. This means that caching unnecessary data can lead to eviction and increased re-computation time.

## References

- [1] Nacéra Seghouani. Massive graph management and analytics. Lecture Notes.