# – Project Report

Submitted in partial fulfilment of the requirements for the award of degree of



**Topic: GraphX Academic Citation Network Analysis**

**B TECH (CSE)**

**Submitted To: Dr. Aryan Maan**

**Submitted By:**

| Name :- | Rahul Mishra |
|---|---|
| Reg No:- | 12220689 |
| Roll No :- | 25 |
| Section :- | K22GB |
| Subject :- | Cluster Computing |
| Course Code :- | INT315 |

**LOVELY PROFESSIONAL UNIVERSITY**

**PHAGWARA, PUNJAB**

# **Acknowledgement**

I would like to express my sincere gratitude to **Dr. Aryan Maan** for valuable guidance and support throughout this project. I am also thankful to Lovely Professional University for providing the facilities and resources to complete this work successfully.

Rahul Mishra

15/11/25

## Abstract

Academic research is interconnected through citations. By treating these citations as a graph, meaningful insights can be extracted — such as identifying influential papers, discovering isolated works, and detecting collaborative clusters. In this project, Apache Spark GraphX is used to build and analyze an academic citation network. The analysis includes in-degree calculation, PageRank scoring, isolated-node detection, and network visualization using Python and NetworkX. A dashboard is also created to summarize the results.

# Table of Contents

## 1. Introduction

Citations define how academic research builds upon previous knowledge. Each citation can be represented as a directional link between two papers. When aggregated, these connections form a *citation network*. Network analysis allows:

- Discovery of influential works

- Identification of emerging research areas

- Understanding of collaboration patterns

- Finding isolated or independent research efforts

Apache Spark's **GraphX** is ideal for processing such large-scale graph datasets due to its distributed nature and built-in graph algorithms like PageRank.

This project uses a real-world citation dataset from AMiner (20,000 edges) and performs an end-to-end analysis and visualization pipeline.

## 2. Problem Statement

Traditional citation analysis tools struggle with:

- Large datasets

- Complex network relationships

- Distributed processing needs

**Goal:** Build a distributed GraphX pipeline that efficiently analyzes an academic citation network.

## 3 . Objectives

1. Load and preprocess citation dataset.

2. Construct directed graph using GraphX.

3. Compute in-degree (citation counts).

4. Run PageRank to find influential papers.

5. Identify isolated papers.

6. Visualize the graph using Python.

7. Build a dashboard illustrating important metrics.

## 4. Technologies Used

| Technology | Purpose |
| --- | --- |
| Apache Spark | Distributed processing |
| GraphX | Graph computation engine |
| Scala | Core backend code |
| sbt | Build tool |
| Python | Visualization |
| NetworkX | Graph drawing |
| Streamlit | Dashboard creation |

# 5. Methodology

## 5.1 Dataset Description

The dataset has been derived from AMiner Citation Network V2 and contains:

- 20,000 directed citation edges
- 6,000 papers

## 5.2 Workflow Steps

- Dataset loaded as RDD
- Graph built using GraphX
- Degrees computed
- PageRank executed
- Isolated nodes detected safely
- Graph exported for visualization
- Python scrip renders the network
- Dashboard summarizes statistics

## 5.3 GraphX Concepts Used

**Vertices**

Each paper is a vertex.

**Edges**

A directed edge exists from A → B if A cites B.

**In-degree**

Number of incoming edges → number of citations received.

**Out-degree**

Number of outgoing edges → number of citations made.

**PageRank**

Measures influence using recursive importance weights.

# 6. Implementation

## 6.1 Project Structure

```
C:\GraphX-Citation\
  ├── build.sbt
  ├── data\
  |    └── citations.txt
  ├── output\
  |    └── graph_edges\
  ├── visualization\
  |    ├── visualize_graph.py
  |    └── dashboard.py
  └── src\
       └── main\
           └── scala\
               └── CitationGraphApp.scala
```

## 6.2 Scala Code

```scala
1   import org.apache.spark.sql.SparkSession
2   import org.apache.spark.graphx._
3   import org.apache.spark.rdd.RDD
4
5   object CitationGraphApp {
6
7     def main(args: Array[String]): Unit = {
8
9       val spark = SparkSession.builder()
10        .appName("GraphX Academic Citation Network")
11        .master("local[*]")
12        .getOrCreate()
13
14      val sc = spark.sparkContext
15
16      // ==============================
17      // 1. LOAD CITATION DATA
18      // ==============================
19      val raw = sc.textFile("data/citations.txt").filter(_.trim.nonEmpty)
20
21      val edges: RDD[Edge[Int]] = raw.map { line =>
22        val parts = line.trim.split("\\s+")
23        Edge(parts(0).toLong, parts(1).toLong, 1)
24      }
25
26      val defaultAttr = 1
27      val graph = Graph.fromEdges(edges, defaultAttr).cache()
28
29      println("=======================================")
30      println(" Graph Loaded Successfully ")
31      println("=======================================")
32      println(s"Total vertices: ${graph.numVertices}")
33      println(s"Total edges:    ${graph.numEdges}")
34      println("=======================================\n")
35
36
37      // ==============================
38      // 2. IN-DEGREE (CITATION COUNT)
39      // ==============================
40      val inDeg = graph.inDegrees.cache()
41
42      println("\nTop 15 Most Cited Papers (In-Degree):")
43      inDeg.sortBy(_._2, ascending = false).take(15).foreach(println)
44
45      println("\n=======================================\n")
46
47
48      // ==============================
49      // 3. PAGE RANK
50      // ==============================
51      println("Running PageRank...")
52      val ranks = graph.pageRank(0.0001).vertices
53
54      val joined = ranks.join(inDeg).map {
55        case (paperId, (rank, indegree)) => (paperId, rank, indegree)
```

```scala
5   object CitationGraphApp {
7     def main(args: Array[String]): Unit = {
53
54      val joined = ranks.join(inDeg).map {
55        case (paperId, (rank, indegree)) => (paperId, rank, indegree)
56      }
57
58      println("\nTop 15 Influential Papers by PageRank:")
59      joined.sortBy(_._2, ascending = false).take(15).foreach(println)
60
61      println("\n=======================================\n")
62
63
64      // ==============================
65      // 4. OUT-DEGREES
66      // ==============================
67      val outDeg = graph.outDegrees.cache()
68
69
70      // ==============================
71      // 5. ISOLATED NODES (SAFE METHOD)
72      // ==============================
73
74      // all vertex IDs
75      val allVerts = graph.vertices.map(_._1)
76
77      // vertices referenced in in- or out- degree
78      val connected = inDeg.map(_._1).union(outDeg.map(_._1)).distinct()
79
80      // isolated = allVerts - connected
81      val isolated = allVerts.subtract(connected)
82
83      println("\nSample Isolated Vertices:")
84      isolated.take(20).foreach(println)
85
86      println("\n=======================================\n")
87
88
89      // ==============================
90      // 6. EXPORT EDGES FOR VISUALIZATION
91      // ==============================
92      graph.edges
93        .map(e => s"${e.srcId},${e.dstId}")
94        .coalesce(1)
95        .saveAsTextFile("output/graph_edges")
96
97      println("Visualization edges saved to: output/graph_edges/")
98      println("Use Gephi / Python NetworkX to visualize.")
99
100     println("\n************ TASKS COMPLETE ************\n")
101
102     spark.stop()
103   }
104 }
105
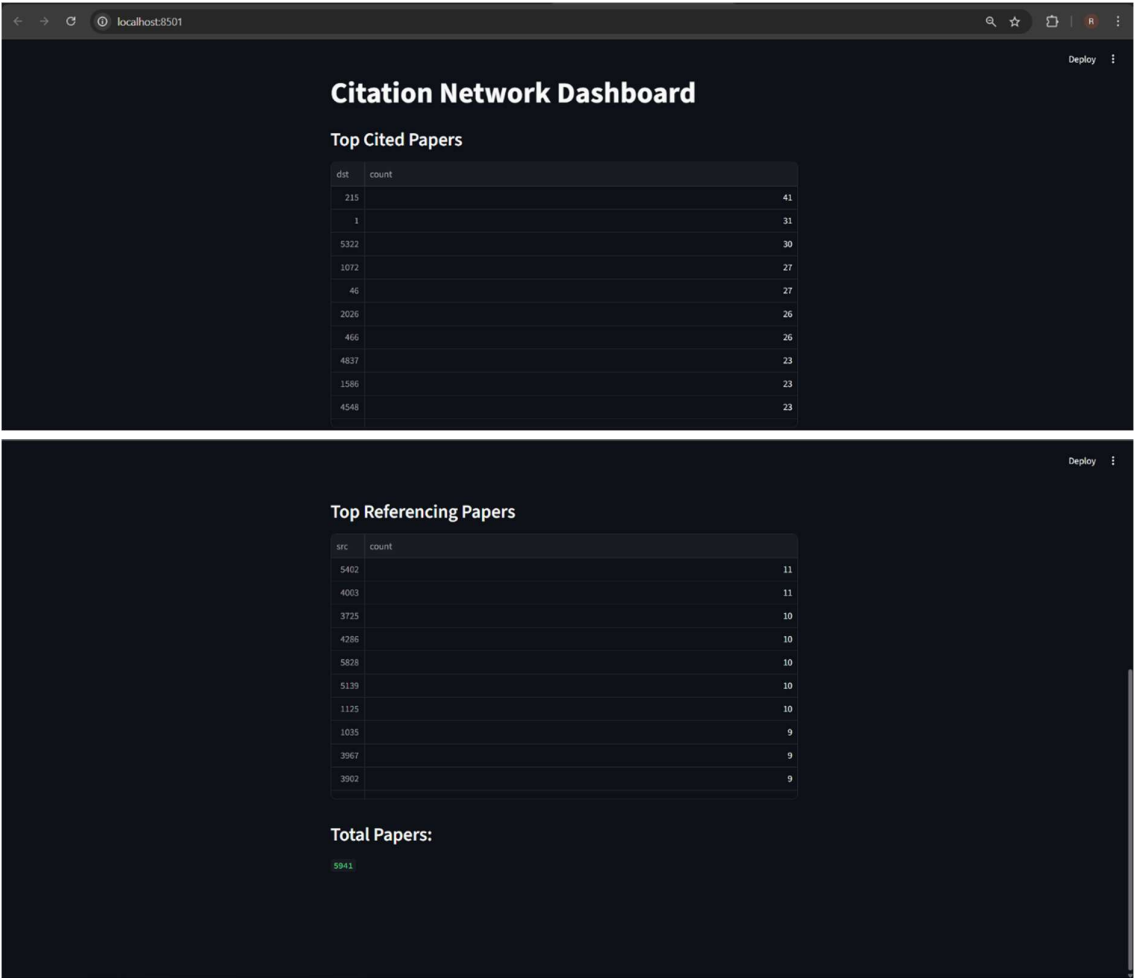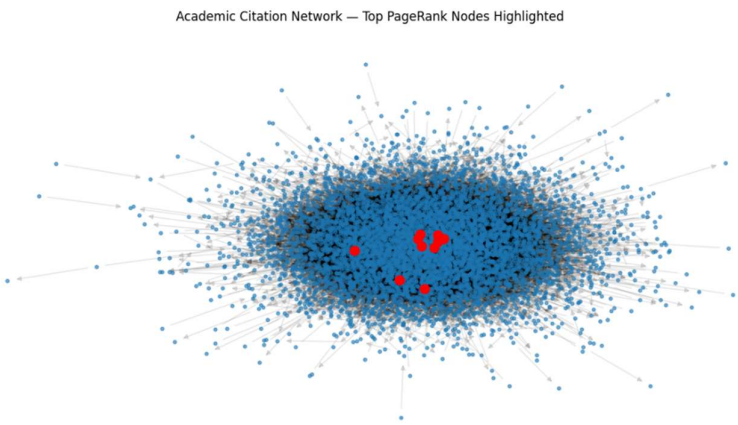```

## 6.3 Python Visualization Script

```python
visualize_graph.py > ...
 1   import networkx as nx
 2   import matplotlib.pyplot as plt
 3
 4   # Load edges
 5   path = r"C:\GraphX-Citation\output\graph_edges\part-00000"
 6
 7   G = nx.DiGraph()
 8
 9   with open(path, "r") as f:
10       for line in f:
11           src, dst = line.strip().split(",")
12           G.add_edge(src, dst)
13
14   # Compute PageRank (small graphs only)
15   pr = nx.pagerank(G)
16   top_nodes = sorted(pr.items(), key=lambda x: x[1], reverse=True)[:10]
17   top_set = {n for (n, _) in top_nodes}
18
19   plt.figure(figsize=(12, 10))
20   pos = nx.spring_layout(G, k=0.25)
21
22   nx.draw_networkx_nodes(G, pos, node_size=10, alpha=0.6)
23   nx.draw_networkx_nodes(G, pos, nodelist=list(top_set), node_color="red", node_size=80)
24   nx.draw_networkx_edges(G, pos, alpha=0.1)
25
26   plt.title("Academic Citation Network — Top PageRank Nodes Highlighted")
27   plt.axis("off")
28   plt.show()
29
```

## 6.4 Dashboard Code

```python
dashboard.py > ...
 1   import streamlit as st        Pin selection to current chat prompt (Ctrl+Alt+X) | Don't show this again (Alt+/)
 2   import pandas as pd
 3
 4   df = pd.read_csv("C:/GraphX-Citation/output/graph_edges/part-00000",
 5                    names=["src","dst"])
 6
 7   st.title("Citation Network Dashboard")
 8
 9   st.subheader("Top Cited Papers")
10   st.write(df["dst"].value_counts().head(20))
11
12   st.subheader("Top Referencing Papers")
13   st.write(df["src"].value_counts().head(20))
14
15   st.subheader("Total Papers:")
16   st.write(len(pd.unique(df[["src","dst"]].values.ravel())))
17
```

Ln 1, Col 17 (9 selected)    Spaces: 4    UTF-8    CRLF    {} Python    3.12.2    Go Live    Prettier

# 7. Results and Screenshots



Academic Citation Network — Top PageRank Nodes Highlighted



**Citation Network Dashboard**

**Top Cited Papers**

| dst | count |
|---|---|
| 215 | 41 |
| 1 | 31 |
| 5322 | 30 |
| 1072 | 27 |
| 46 | 27 |
| 2026 | 26 |
| 466 | 26 |
| 4837 | 23 |
| 1586 | 23 |
| 4548 | 23 |

**Top Referencing Papers**

| src | count |
|---|---|
| 5402 | 11 |
| 4003 | 11 |
| 3725 | 10 |
| 4286 | 10 |
| 5828 | 10 |
| 5139 | 10 |
| 1125 | 10 |
| 1035 | 9 |
| 3967 | 9 |
| 3902 | 9 |

**Total Papers:**

5941

## 8. Conclusion

This project successfully demonstrates:

- How Spark GraphX can be used to analyze citation networks
- Performing PageRank on distributed graph data
- Computing academic influence metrics
- Identifying isolated or unique research works
- Visualizing complex networks using Python tools
  The end-to-end pipeline is efficient, scalable, and suitable for large datasets.