

Analyzing Amazon's Co-Purchasing Network With Apache Spark

Dipesh Mishra
CSE
IIIT Dharwad, India
20bcs043@iiitdwd.ac.in

Mahendra Singh Puniya
CSE
IIIT Dharwad, India
20bcs082@iiitdwd.ac.in

Ravikant
CSE
IIIT Dharwad, India
20bcs110@iiitdwd.ac.in

Rishabh Gautam
CSE
IIIT Dharwad, India
20bcs112@iiitdwd.ac.in

Abstract—This analysis delves into the Amazon Product Co-purchasing Network, a dataset collected in June 2003, to unravel patterns of customer behavior and product relationships on the Amazon platform. Leveraging Apache Spark and the Connected Components algorithm, this study identifies clusters of interconnected products, shedding light on co-purchasing behavior. The findings have profound implications for enhancing recommendation systems and optimizing marketing strategies. The report underscores the significance of Apache Spark in efficiently processing and analyzing large-scale network data, offering insights that empower data-driven decision making for e-commerce applications.

Keywords—Amazon Product Network, Network Analysis, Apache Spark, Connected Components, Customer Behavior, Product Relationships, Recommendation Systems, Marketing Strategies, E-commerce Data, Data Processing, Insights, Graph Analysis, Network Statistics

I. INTRODUCTION

In the landscape of big data and distributed computing, Apache Spark has emerged as a powerful framework for processing and analyzing vast datasets efficiently. This report centers on the analysis of the Amazon Product Co-purchasing Network using Apache Spark, an open-source, lightning-fast, and cluster computing system. The dataset, collected in June 2003, was harvested by crawling the Amazon website and is rooted in the "Customers Who Bought This Item Also Bought" feature, a cornerstone of e-commerce recommendation systems. The Amazon Product Co-purchasing Network provides a wealth of insights into customer interactions with products on the Amazon platform. By leveraging Apache Spark, we unlock the capability to process and analyze this vast and interconnected dataset with ease. Our analysis, powered by the Connected Components algorithm in Apache Spark's GraphX library, allows us to uncover patterns and relationships within the network, enabling us to understand how products are co-purchased, clustered, recommended together.

II. DATASET INFORMATION

The dataset under consideration is the Amazon Product Co-purchasing Network, collected in June 2003. This network was obtained through web crawling of the Amazon website, relying on the "Customers Who Bought This Item Also Bought" feature. In this network, if a product 'i' is frequently co-purchased with product 'j', a directed edge is

created from 'i' to 'j'. The dataset provides essential statistics, offering valuable insights into its scale and structure:

- **Nodes: 403,394**
- **Edges: 3,387,388**
- **Average Clustering Coefficient: 0.4177**
- **Number of Triangles: 3,986,507**
- **Fraction of Closed Triangles: 0.06206**
- **Diameter (Longest Shortest Path): 21**
- **90-Percentile Effective Diameter: 7.6**

TABLE I.

DATASET	Nodes	Edges
<i>Largest Weakly Connected Component (WCC)</i>	403,364 (100% of nodes)	3,387,224 (100% of edges)
<i>Largest Strongly Connected Component (SCC)</i>	395,234 (98% of nodes)	3,301,092 (97.5% of edges)

These statistics provide a comprehensive overview of the dataset's scale, connectivity, and inherent structure, which is essential for conducting meaningful network analysis and drawing actionable insights from the data.

III. WORKFLOW

A. Data Acquisition

Obtain the Amazon Product Co-purchasing Network dataset, collected in June 2003, either from the source provided or the relevant data repository.

B. Environment Setup

Set up your development environment with Apache Spark and relevant libraries, ensuring that your infrastructure can handle the size and complexity of the dataset.

C. Data Loading

Load the dataset into your Spark environment using the GraphX library. Ensure that the data is appropriately structured for analysis.

D. Connected Components Analysis

Utilize the Connected Components algorithm from GraphX to identify connected components within the network. This step will label each component with the ID of its lowest-numbered vertex.

