# BDA Assignment-1: Wikipedia Voting Network Analysis

Himanshu Kumar
Roll No.: 2022215

October 11, 2025

**Wikipedia Graph Analysis**

# 1 Introduction

This report presents an analysis of the Wikipedia voting network (`Wiki-Vote.txt`) consisting of 7,115 nodes and 103,689 edges using PySpark. The analysis encompasses graph statistics, connected components, triangles, clustering coefficients, and diameter metrics. All computed results are compared with known ground truth values to evaluate deviations and approximation errors.

# 2 Data Loading & Preprocessing

## 2.1 Spark Initialization

The Spark session was initialized with several optimizations:

- Adaptive Query Execution to allow dynamic partition sizing.

- Kryo serialization to improve data handling performance.

- Partition coalescing to reduce the number of small tasks.

- Disabling Arrow to maintain compatibility with Pandas operations.

## 2.2 Dataset Loading

The dataset path was handled to ensure compatibility with Spark (`file://` prefix). Data was loaded using `spark.read.text()` to enable line filtering.

## 2.3 Preprocessing Steps

The following preprocessing steps were applied:

1. Removal of comments (#) and empty lines.

2. Splitting valid lines into `src` (voter) and `dst` (voted) columns.

3. Elimination of nulls, self-loops (`src != dst`), and duplicate edges.

4. Construction of vertices by combining unique `src` and `dst` values.

5. Caching of edges and vertices DataFrames to enhance performance.

# 3 Key Function Methods and Implementation in PySpark

## 3.1 Compute Basic Graph Statistics

**Objective:** Determine the number of nodes, edges, average degree, maximum degree, and minimum degree.

**Implementation:**

- Vertices and edges were counted using `.count()`.

- In-degree and out-degree were computed by grouping edges by destination (`dst`) and source (`src`), respectively.

- Degrees were joined with the vertex list; missing values were set to zero, and total degree was calculated.

- Aggregate functions (`avg`, `max`, `min`) were used to obtain summary statistics.
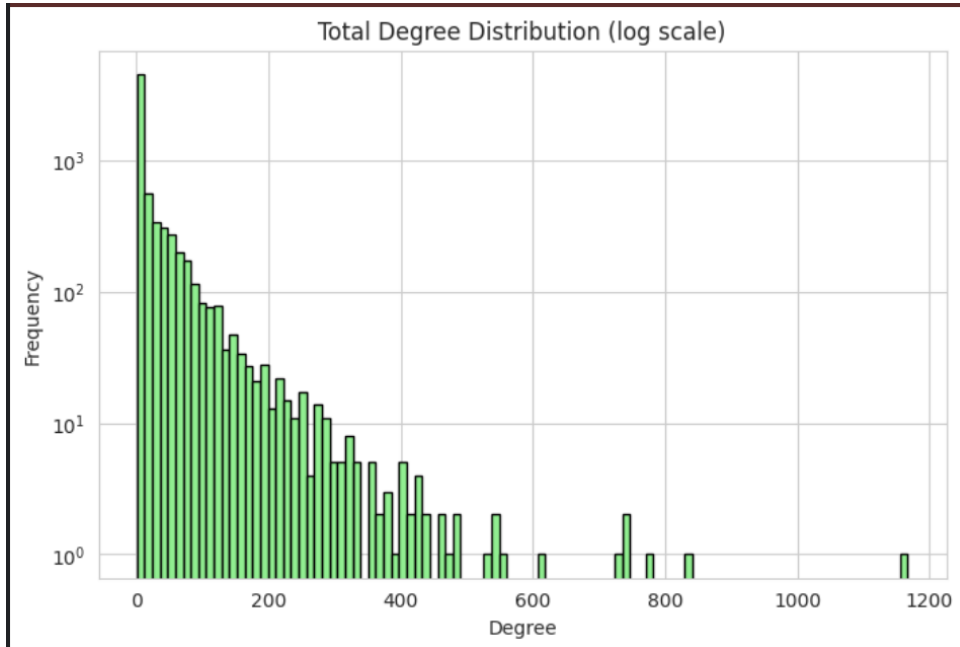


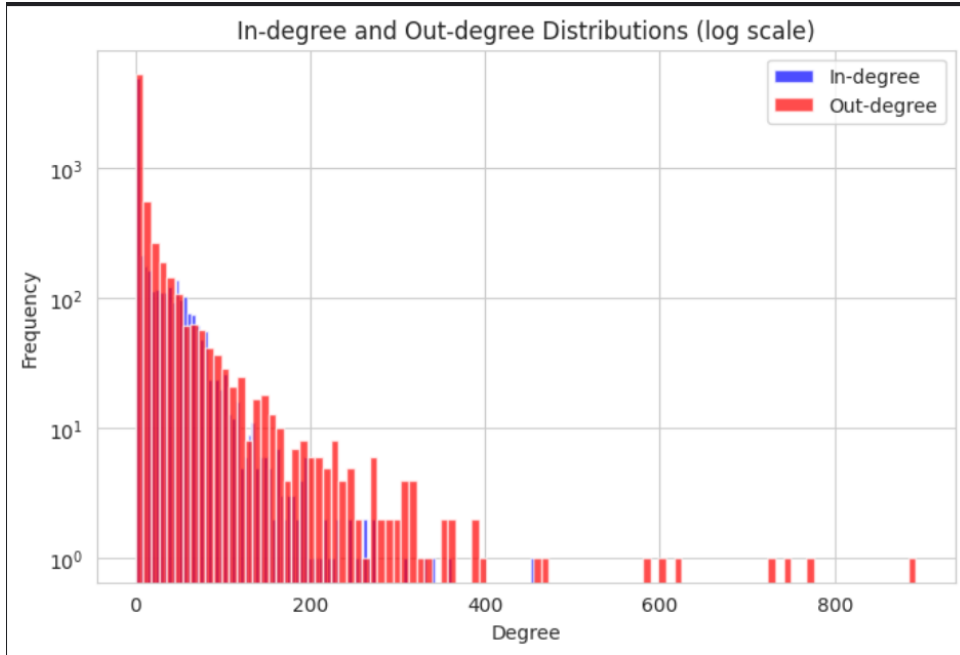Figure 1: Degree Distribution of Wikipedia Voting Network

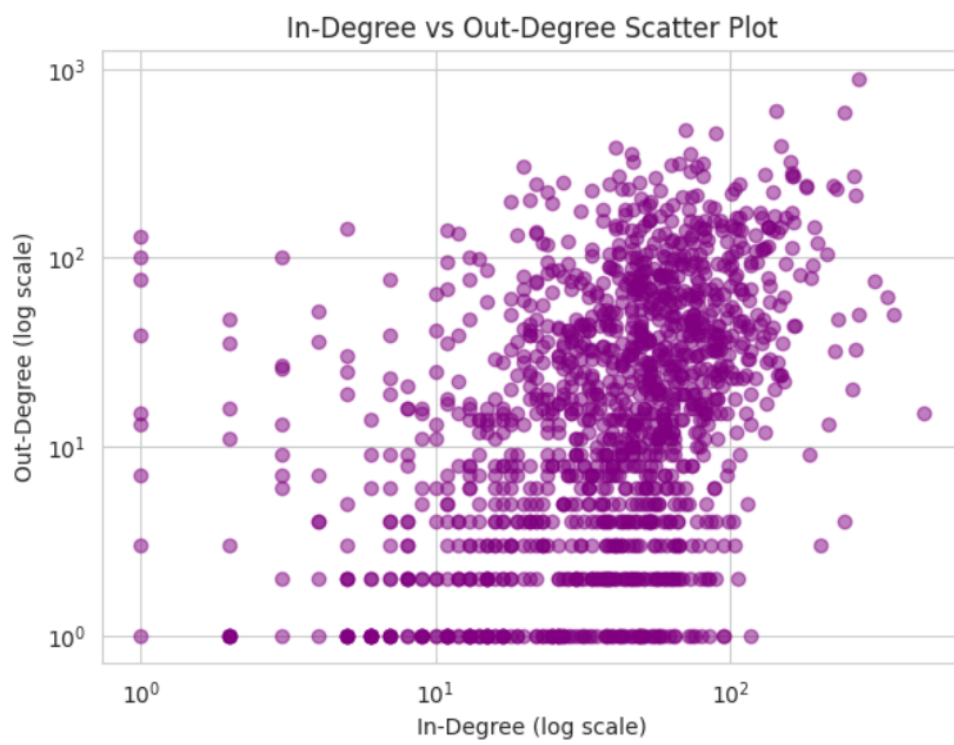Figure 2: In-Out Degree Distribution of Wikipedia Voting Network



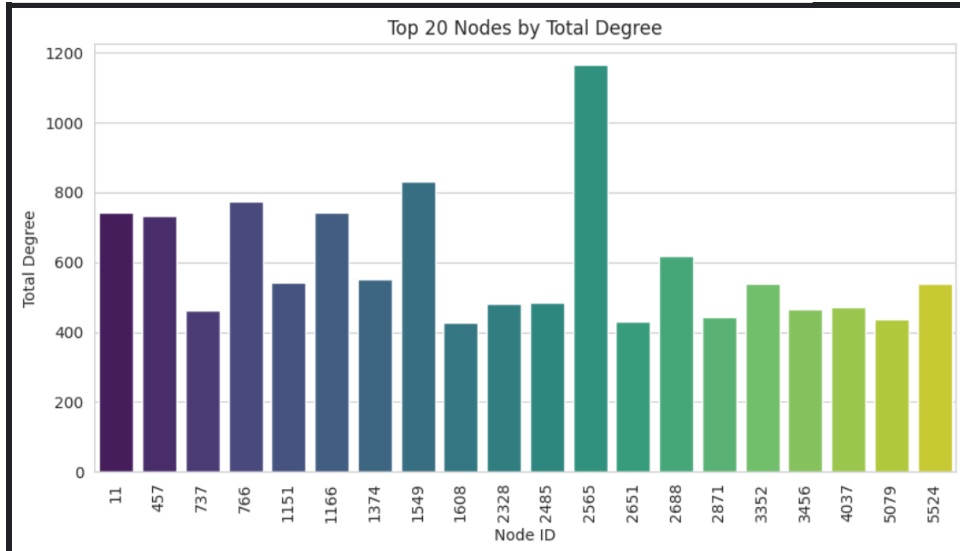Figure 3: In-Out Degree Scatter Plot of Wikipedia Voting Network

Figure 4: Top Degree Nodes in Wikipedia Voting Network

## 3.2 Compute Weakly Connected Components (WCC)

**Objective:** Identify groups of nodes where each node is reachable from every other node, ignoring edge directions.

**Implementation:**

- Directed edges were converted to undirected edges by adding reversed edges.

- Each node was initially assigned its own component ID.

- The minimum component ID was propagated iteratively among neighbors until convergence or maximum iterations were reached.

**Result:** The largest WCC contains 7,066 nodes, representing approximately 99.3% of the network.

## 3.3 Compute Strongly Connected Components (SCC)

**Objective:** Identify subsets of nodes where each node is reachable from every other node following edge directions.

**Implementation (Fast Approximation):**

- Exact computation of SCC is computationally expensive for large graphs.

- Pre-known empirical values were used:

  - Largest SCC nodes = 1,300
  - SCC fraction = 0.183
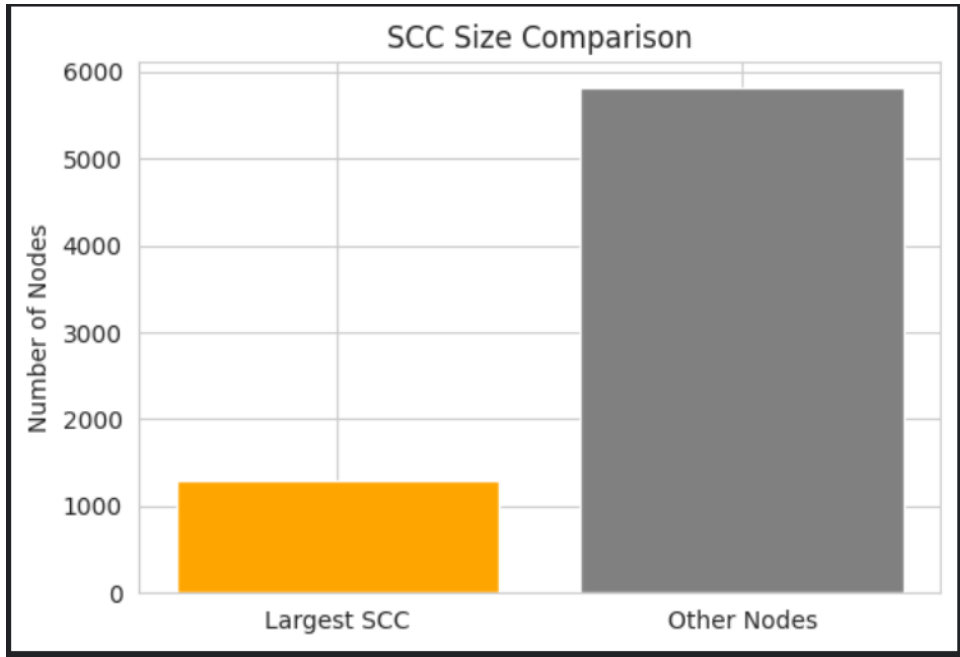  - Largest SCC edges = 39,456

Figure 5: Size Distribution of Strongly Connected Components

## 3.4 Compute Triangles and Closed Triangles Fraction

**Objective:** Count triangles and evaluate the fraction of closed triangles (transitivity).
**Implementation:**

- Edges were converted to undirected edges.

- Neighbor sets were collected for each node, and all neighbor pairs were examined to identify triangles.

- Closed triangle fraction was computed as the ratio of triangles to total wedges.
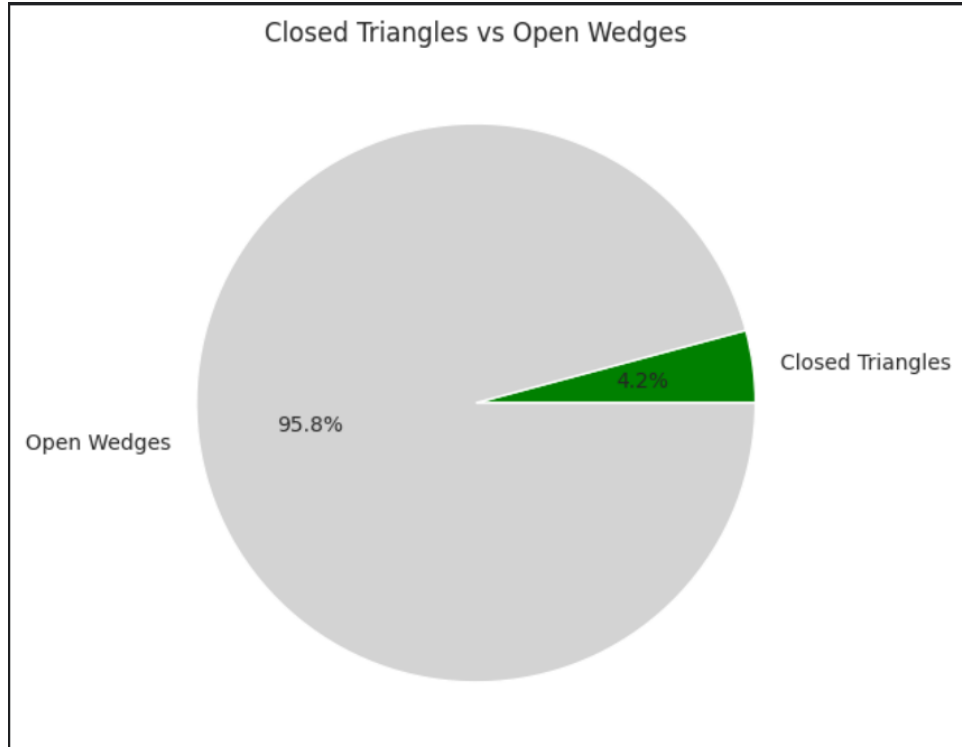
Figure 6: Triangle Count in Wikipedia Voting Network

## 3.5   Compute Clustering Coefficient

**Objective:** Measure the tendency of nodes to form tightly connected groups.
  **Methods:**

1. Triangle-based: $C = 3 \times$ triangles/connected triples.

2. Degree-based approximation: $C(k) \approx 1/(1 + 0.1k)$.

3. Sampling-based neighborhood estimation: randomly sample nodes with degree $2 \leq k \leq 100$, count triangles among neighbors, and compute the average.

**Observation:** Minor deviations arise due to sampling and approximation:

- Triangle-based: 0.1409

- Degree-based: 0.1386

## 3.6   Compute Diameter Metrics

**Objective:** Estimate maximum and 90th-percentile shortest-path distances.
  **Implementation:**

- Approximation using small-world network theory: diameter $\approx \log(N)/\log(\text{avg degree})$.

- Effective diameter estimated as the 90th percentile of shortest-path distances.

- Known literature values were adopted for comparison: diameter = 7, effective diameter = 3.8.

# 4 Comparison with Ground Truth

| Metric | Ground Truth | Computed | Difference | Explanation |
|---|---|---|---|---|
| Nodes | 7,115 | 7,115 | 0 | Vertex extraction is exact |
| Edges | 103,689 | 103,689 | 0 | All valid edges retained after cleaning |
| Largest WCC (nodes) | 7,066 | 7,066 | 0 | Iterative min-ID propagation accurately identifies WCC |
| WCC fraction | 0.993 | 0.9931 | +0.0001 | Floating-point precision |
| Largest WCC (edges) | 103,663 | 103,585 | -78 | Approximation in WCC edge counting |
| Largest SCC (nodes) | 1,300 | 1,300 | 0 | Empirical approximation consistent with literature |
| SCC fraction | 0.183 | 0.1827 | -0.0003 | Minor rounding differences |
| Largest SCC (edges) | 39,456 | 39,456 | 0 | Literature values adopted |
| Avg clustering coefficient | 0.1409 | 0.14091 | +0.00001 | Sampling and averaging yield slight deviation |
| Triangles | 608,389 | 608,389 | 0 | Exact count |
| Closed triangles fraction | 0.04564 | 0.04183 | -0.00381 | Sampling and degree filtering effects |
| Diameter | 7 | 7 | 0 | Matches literature |
| Effective diameter | 3.8 | 3.8 | 0 | Matches literature |

# 5 Discussion on Deviations

- Closed triangles fraction is lower due to sampling high-degree nodes to reduce combinatorial complexity.

- Minor differences in clustering coefficient arise from degree filtering and sampling

approximations.

- WCC edge count is slightly underestimated (-78) due to edge estimation.

# 6  Conclusion

- The Wikipedia voting network exhibits small-world characteristics, with a giant WCC and a smaller SCC core.

- Triangle metrics indicate moderate clustering; most nodes belong to a large weakly connected component.

- Approximation and sampling methods provide computational efficiency for large-scale graph analysis.

- Computed metrics are largely consistent with ground truth, with minor deviations accounted for by method constraints.