

Generation of Test Datasets for Memory and Performance Analysis

For performance and memory usage analysis, I generated test datasets of different sizes using custom Python scripts. These datasets simulate networks of professors and the courses they teach, creating a realistic graph structure for testing the ProfessorPaths program.

Test Datasets

I created two primary test files:

- `data/test_medium.csv` - Medium-sized network with 500 professors and 100 courses
- `data/large_network.csv` - Large network with 5,000 professors and 1,000 courses

Generation Method

1. I created two Python scripts to generate the datasets:
 - `data/generate_medium.py` - For the medium test dataset
 - `data/generate_larger.py` - For the large test dataset
2. The scripts create random assignments of professors to courses with realistic distribution patterns:
 - Each professor teaches between 3-8 courses (realistic teaching load)
 - Course popularity follows a power-law distribution (some courses have many professors)
 - The data is written in CSV format with professor-course pairs
3. Medium test characteristics:
 - 500 professors
 - 100 unique courses
 - Approximately 5,000 professor-course assignments
4. Large test characteristics:
 - 5,000 professors
 - 1,000 unique courses
 - Approximately 25,000 professor-course assignments
 - Creates approximately 1.5 million potential edges in the graph

Performance Analysis Results

I tested the ProfessorPaths implementation with these datasets and observed the following results:

Medium Test (500 professors)

- Graph creation time: ~55 milliseconds
- Memory usage: ~20 MB
- Path finding operations: Under 50 milliseconds on average

Large Test (5,000 professors)

- Graph creation time: 308 milliseconds
- Memory usage: 178 MB (within 256 MB limit)
- Path finding times: 150-410 milliseconds per query
- All paths found successfully within memory constraints

Optimizations

To improve performance and reduce memory usage, I implemented several optimizations:

1. **Disabled checkRep() for Performance Testing:** The `CHECK_REP_ENABLED` flag in `Graph.java` was set to `false` to eliminate the expensive representation invariant checks during performance-critical operations.
2. **BFS Optimization:** Improved the breadth-first search algorithm by separating professor node sorting from course sorting, ensuring lexicographical ordering while minimizing performance impact.
3. **Memory-Efficient Data Structures:** Used efficient collections for path tracking and queue management during BFS traversal.

Conclusion

The implementation successfully processes the large network dataset (5,000 professors) within the 256 MB memory constraint and finds paths in "a fraction of a second" as required. The performance metrics demonstrate that the chosen data structures and algorithms are efficient even with large datasets, validating the design decisions made in the implementation.