# Using Big Data to Connect the Dots from One Place to Another

## Chuck Cartledge

## June 21, 2016 at 1:15pm

## Table of contents I

| Historical origins | Neo4j | Results | Explorations | Playing | Conclusion | References |
|---|---|---|---|---|---|---|
| ●0000 | 0000000 | 00 00000 | 00000 | 00 | | |

Terms and definitions

## Königsberg and Euler [1]

- Popular Sunday puzzle was to cross all bridges exactly once
- Mayor asks Leonard Euler to solve the puzzle
- Euler declines, thought the problem was trivial
- Euler changes mind.
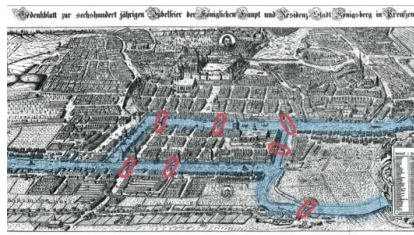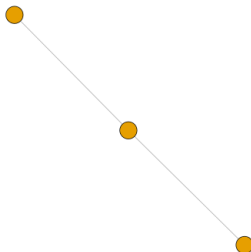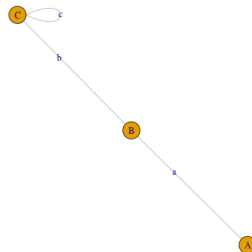- Euler publishes paper in 1736 detailing impossible solution, formulates general solution

Birth place of graph theory.



Image from [1].

Historical origins | Neo4j | Results | Explorations | Playing | Conclusion | References
○●○○○ | ○○○○○○○ | ○○ | ○○○○○ | ○○ | |
| | ○○○○○ | | |

Terms and definitions

# Types of graphs (1 of 2)



(a) Simple

(b) Labeled

## Types of graphs (2 of 2)



(c) Directed



(d) Useful

Historical origins
○○○●○

Neo4j
○○○○○○○

Results
○○
○○○○○

Explorations
○○○○○

Playing
○○

Conclusion
○○

References

Terms and definitions

## ODU Computer Science undergraduate prerequisites

Historical origins
○○○○●

Neo4j
○○○○○○○

Results
○○
○○○○○

Explorations
○○○○○

Playing
○○

Conclusion

References

Terms and definitions

## Terms and definitions

- nodes ≡ vertices
- arcs ≡ edges
- arcs ≡ relationships

Terms used interchangeably.

Neo4j uses the terms nodes and relationships.

| Historical origins | Neo4j | Results | Explorations | Playing | Conclusion | References |
| ----- | ----- | ----- | ----- | ----- | ----- | ----- |
| ○○○○○ | ●○○○○○○ | ○○ | ○○○○○ | ○○ | | |
| | | ○○ | ○○○○○ | | | |

The basics.

# Where to get it.



Figure: Available with different licenses from http://neo4j.com/

Historical origins
○○○○○

Neo4j
○●○○○○○○

Results
○○
○○○○○

Explorations
○○○○○

Playing
○○

Conclusion

References

The basics.

## How to start it.



Figure: Starting neo4j manually.

Historical origins
○○○○○

Neo4j
○○○●○○○○

Results
○○
○○○○○

Explorations
○○○○○

Playing
○○

Conclusion

References

The basics.

# How to interact with it (1 of 4).



Figure: Via a browser. Note the URL.

| Historical origins | Neo4j | Results | Explorations | Playing | Conclusion | References |
| ----- | ----- | ----- | ----- | ----- | ----- | ----- |
| 00000 | 0000●000 | 00 | 00000 | 00 | | |
| | | 00000 | | | | |

The basics.

How to interact with it (2 of 4).



Figure: Information about the database.

| Historical origins | Neo4j | Results | Explorations | Playing | Conclusion | References |
|---|---|---|---|---|---|---|
| ○○○○○ | ○○○○●○○ | ○○ | ○○○○○ | ○○ | ○○ | |
| | | ○○○○○ | | | | |

The basics.

How to interact with it (3 of 4).



Figure: A command line interface.

# How to interact with it (4 of 4).



Figure: Help on the CLI.

| Historical origins | Neo4j | Results | Explorations | Playing | Conclusion | References |
|--------------------|-------|---------|--------------|---------|------------|-----------|
| ○○○○○ | ○○○○○○● | ○○ | ○○○○○ | ○○ | | |
| | | ○○○○○ | | | | |

The basics.

## Neo4j is idempotent

Neo4j nodes can be thought of as a set where each set member is unique/citemathDictionary.

- The time to insert an unconstrained node can be: $O(n^2)$.

- The time to insert a constrained node can be: $O(n)$.

- The time to execute a query can be: $O(x)$

- The time to re-execute a previous query can be: $O(c)$ (for a small $c$)

Take away: constrain your nodes.



Image from [2].

| Historical origins | Neo4j | **Results** | Explorations | Playing | Conclusion | References |
| :--- | :--- | :--- | :--- | :--- | :--- | :--- |
| ○○○○○ | ○○○○○○○ | ●○ | ○○○○○ | ○○ | | |
| | | ○○○○○ | | | | |

Failure

## IMDb (1 of 2)

Table: Time spent loading the IMDb via the neo4j-shell. The time to load the entire IMDb database was too long to be practical.

| Size | % IMDb | R-seconds | Neo4j seconds | Nodes loaded |
| ---: | ---: | ---: | ---: | ---: |
| 8,000 | - | 14.070 | 2.6 | 2,300 |
| 80,000 | - | 57.462 | 8.820 | 26,041 |
| 800,000 | 4 | 609.890 | 49.250 | 198,998 |
| 8,000,000 | 45 | 1,245.048 | 556.452 | 649,309 |
| 80,000,000 | 100 | 43,600,386.000 | 1,840.666 | 1,466,720 |

# IMDB (2 of 2)



Figure: JVM needs swap space.

# OpenFlight.org



Figure: OpenFlight.org homepage. http://openflights.org/

| Historical origins | Neo4j | **Results** | Explorations | Playing | Conclusion | References |
|---|---|---|---|---|---|---|
| ○○○○○ | ○○○○○○○ | ○○ | ○○○○○ | ○○ | ○○ | |
| | | ○●○○○○ | | | | |

Success

OpenFlight.org database.



Figure: Interested in airport locations, and service between airports.
http://openflights.org/data.html

| Historical origins | Neo4j | **Results** | Explorations | Playing | Conclusion | References |
| :-- | :-- | :-- | :-- | :-- | :-- | :-- |
| ○○○○○ | ○○○○○○○ | ○○ | ○○○○○ | ○○ | ○○ | |
| | | ○○○●○○ | | | | |

Success

Sample airport data.



Figure: Sample data. There are 8,109 airports identified.

| Historical origins | Neo4j | Results | Explorations | Playing | Conclusion | References |
| 00000 | 0000000 | 00 | 00000 | 00 | | |
| | | 000●0 | | | | |

Success

Sample route data.



Figure: Sample route data. There are 67,665 route records. Not all airports have service.

Airports around the world.

| Historical origins | Neo4j | Results | Explorations | Playing | Conclusion | References |
| 00000 | 0000000 | 00 | ●0000 | 00 | | |
| | | 00000 | | | | |

Available data

## Choosing a source and terminating airport.

Available airport data:

- AirportID: Unique OpenFlights identifier for this airport.
- Name: Name of airport. May or may not contain the City name.
- City: Main city served by airport. May be spelled differently from Name.
- Country: Country or territory where airport is located.
- IATA: 3-letter FAA code, for airports located in Country "United States of America". 3-letter IATA code, for all other airports. Blank if not assigned.
- ICAO: 4-letter ICAO code. Blank if not assigned.
- Latitude: Decimal degrees
- Longitude: Decimal degrees
- Altitude: In feet.
- Timezone: Hours offset from UTC
- DST: Daylight savings time
- Timezone: Timezone name

| Historical origins | Neo4j | Results | Explorations | Playing | Conclusion | References |
| ----- | ----- | ----- | ----- | ----- | ----- | ----- |
| 00000 | 0000000 | 00 | 0●000 | 00 | | |
| | | 00000 | | | | |

Available data

## Choosing a source and terminating IATA airport.

Steps to identify the International Air Transport Association (IATA) code for exploration:

1. Search the airport location database for airports of interest.
2. Extract the IATA codes for those airports.
3. Give those codes to the airlinePaths.R program.
4. Evaluate results. Results include:
   4.a All airports from the source to the terminating airport.
   4.b The path from source to terminating on four different scaled maps.
   4.c Automatic inclusion of results into the next printing of the final report.

ICAO stands for International Civil Aviation Organization, a UN organization.[3]

Sample tabular data from EYW-YXY (Key West Intl to Whitehorse Intl).

EYW-YXY (Key West Intl to Whitehorse Intl) airports along the flight path.

| Name | IATA | IACO | Lat. | Lon. |
|---|---|---|---|---|
| Key West Intl | EYW | KEYW | 24.556 | −81.760 |
| Orlando Intl | MCO | KMCO | 28.429 | −81.309 |
| Mc Carran Intl | LAS | KLAS | 36.080 | −115.152 |
| Vancouver Intl | YVR | CYVR | 49.194 | −123.184 |
| Whitehorse Intl | YXY | CYXY | 60.710 | −135.067 |

Figure: Tabular data from EYW-YXY (Key West Intl to Whitehorse Intl).

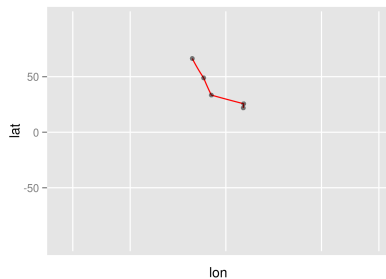| Historical origins | Neo4j | Results | **Explorations** | Playing | Conclusion | References |
| 00000 | 0000000 | 00 | 000●0 | 00 | | |
| | | 00000 | | | | |

Available data

Sample map data from EYW-YXY (1 of 2).



Figure: World view.
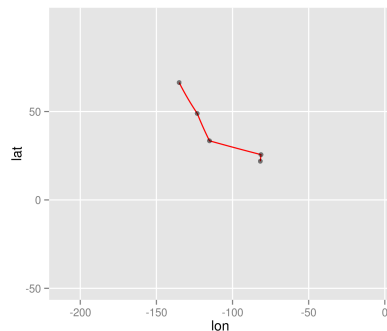


Figure: Hemisphere view.

Sample data from EYW-YXY (2 of 2).



Figure: Continental view



Figure: "Local"

Places that were fun to think about. (1 of 2)

| Historical origins | Neo4j | Results | Explorations | Playing | Conclusion | References |
|---|---|---|---|---|---|---|
| ○○○○○ | ○○○○○○○ | ○○ | ○○○○○ | ○● | | |
| | | ○○○○○ | | | | |

Misc. maps.

# Places that were fun to think about. (2 of 2)

What have we covered?

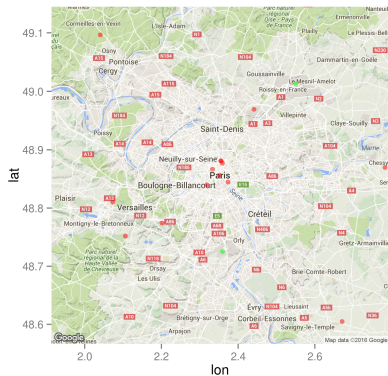- Neo4j is a graph database management system optimized to answer questions that can be framed as a graph.
- Neo4j is Java based, with all the strengths and limitations inherent with a Java Virtual Machine.
- The neo4j community edition has practical limitations on the size of the databases it can support.
- Neo4j is very fast and efficient at answering specific types of questions.



Next time: who knows? Columnar databases?

References I

[1] Teo Paoletti, Leonard euler's solution to the königsberg bridge problem,
    Loci **3** (2011).

[2] BBC Staff, Maths, basic skills, http://www.bbc.co.uk/schools/
    gcsebitesize/maths/algebra/graphsrev5.shtml, 2016.

[3] IACO Staff, About icao,
    http://www.icao.int/about-icao/Pages/default.aspx, 2016.