

# COMS W4111: Introduction to Databases

## Spring 2024, Sections 002/V02

### Homework 4

## Introduction

- This notebook contains HW4. **Both Programming and Nonprogramming tracks should complete this homework.**
- You will submit **PDF and ZIP files** for this assignment. Gradescope will have two separate assignments for these.
- For the PDF:
  - The most reliable way to save as PDF is to go to your browser's menu bar and click `File -> Print`. Switch the orientation to landscape mode, and hit save.
  - **MAKE SURE ALL YOUR WORK (CODE AND SCREENSHOTS) IS VISIBLE ON THE PDF. YOU WILL NOT GET CREDIT IF ANYTHING IS CUT OFF.** Reach out for troubleshooting.
  - **MAKE SURE YOU DON'T SUBMIT A SINGLE PAGE PDF.** Your PDF should have multiple pages.
- For the ZIP:
  - Zip a folder containing this notebook and any screenshots.
  - You may delete any unnecessary files, such as caches.

## Setup

```
In [ ]: %load_ext sql
        %sql mysql+pymysql://root:tzy123456@localhost
```

```
In [ ]: # import sys

        # !{sys.executable} -m pip install --upgrade pymongo
        # !{sys.executable} -m pip install --upgrade neo4j
```

- If you get warnings below, try restarting your kernel

```
In [ ]: import neo4j
        import pandas
        import pymongo

        # TODO: Fill in with your Mongo URL
        mongo_url = "mongodb+srv://dbuser:tzy123456@w4111-hw4.vdub5km.mongodb.net/?retryWrites=true&w=majority&appName=w4111-hw4"
        mongo_client = pymongo.MongoClient(mongo_url)

        # TODO: Fill in with your Neo4j credentials
        neo4j_url = "neo4j+s://cbde716d.databases.neo4j.io"
        neo4j_password = "B--3xRLF6-StPdJ6qX4P3VnmvPF9oimimsUWW3uSw_s"
        # username is always "neo4j"
        graph = neo4j.GraphDatabase.driver(neo4j_url, auth=("neo4j", neo4j_password))
        graph.verify_connectivity()
```

# Written Questions

- As usual, do not bloviate

## W1

Explain the following concepts:

1. Clustering index
2. Nonclustering index
3. Sparse index
4. Dense index

### Answer

1. Clustering index: the index whose search key specifies the sequential order of the file.
2. Nonclustering index: an index whose search key specifies an order different from the sequential order of the file
3. Sparse index: contains index records for only some search-key values.
4. Dense index: Index record appears for every search-key value in the file.

## W2

Explain why nonclustering indexes must be dense.

### Answer

Nonclustering indexes do not store the same order in the file, and it will be hard to find the target record if using sparse index. Sparse index first locates the nearest pointer, and then find the record by scanning sequentially.

## W3

Suppose that, in a table containing information about Columbia classes, the columns `class_code` , `semester` , and `year` are queried frequently **individually**. Would putting a composite index on `(class_code, semester, year)` be a good idea? Why or why not?

### Answer

This is not a good idea. A composite index will be difficult to search for semester and year individually.

## W4

Explain the following concepts:

1. Hash index
2. B+ tree index

### Answer

1. Hash index: Using hashing to store the index. In a hash index, buckets store entries with pointers to records and buckets are derived with hash function.
2. B+ tree index: Using B+ tree to store the index. A B+ tree is a rooted tree that all paths from root to leaf are of the same length and each node that is not a root or a leaf has between  $\lceil n/2 \rceil$  and  $n$  children.

## W5

Give one advantage and one disadvantage of hash indexes compared to B+ tree indexes.

#### Answer

1. Advantage: If hash function being appropriately chosen, search operation can be done in  $O(1)$ .
2. Disadvantage: does not support range query.

## W6

Explain the role of the buffer in a DBMS. Why doesn't the DBMS simply load the entire database in its buffer?

#### Answer

Portion of main memory available to store copies of disk blocks to make operations faster.

Not loading the entire database in its buffer is because main memory is limited.

## W7

Explain the following concepts as they relate to buffer replacement policies:

1. Clairvoyant algorithm
2. Least recently used strategy
3. Most recently used strategy
4. Clock algorithm

#### Answer

1. Clairvoyant algorithm: replaces the blocks that will not be used for the longest period of time in the future. However not implementable due to requirement of future knowledge.
2. Least recently used strategy: tracks the usage of blocks in memory and removes the one that has been unused for the longest time.
3. Most recently used strategy: discards the most recently used page when it needs to free up space in the buffer.
4. Clock algorithm:
  - Sweep second hand clockwise one frame at a time.
  - If bit is 0, choose for replacement.
  - If bit is 1, set bit to zero and go to next frame. Set bit to 1 if being used by application.

## W8

NoSQL databases have become increasingly popular for applications. List 3 benefits of using NoSQL databases over SQL ones.

#### Answer

1. designed to excel in scalability, especially when it comes to big data across different servers.
2. more flexible data model, which is beneficial for applications that handle a variety of data types or rapidly changing schemas.
3. improved performance for certain types of queries due to their non-relational nature.

## W9

Explain the concept between impedance mismatch and how it relates to SQL vs. NoSQL databases.

#### Answer

Impedance mismatch describes the difficulties that arise when trying to integrate object-oriented programming models with relational SQL database models.

NoSQL databases, by contrast, often provide a more natural fit with object-oriented programming due to their flexible schemas.

## W10

The relationship between students and courses is many-to-many. Due to its emphasis on atomicity, modeling this relationship in a relational database would require an associative entity. Explain how this relationship could be modeled in

1. A document database, such as MongoDB
2. A graph database, such as Neo4j

### Answer

1. document database: create a separate collection that explicitly stores the relationships between students and courses.
2. graph database: edges support many-to-many relationship.

---

## MongoDB

- The cell below creates a database `w4111`, then a collection `episodes` inside `w4111`. It then inserts GoT episode data into the collection.

```
In [ ]: import json

with open("episodes.json") as f:
    data = json.load(f)

episodes = mongo_client["w4111"]["episodes"]
episodes.drop()
episodes.insert_many(data)
print("Successfully inserted episode data")
```

Successfully inserted episode data

## M1

- Write and execute a query that shows episodes and the number of scenes they contain
- Your aggregation should have the following attributes:
  - `episodeTitle`
  - `seasonNum`
  - `episodeNum`
  - `numScenes`, which is the length of the episode's `scenes` array
- Order your output on `numScenes` descending, and only keep episodes with more than 100 scenes

```
In [ ]: res = episodes.aggregate([
    {
        '$project': {
            'episodeTitle': 1,
            'seasonNum': 1,
            'episodeNum': 1,
            'numScenes': {
                '$size': '$scenes'
            },
            '_id': 0
        }
    })
```

```
    }, {
      '$match': {
        'numScenes': {
          '$gt': 100
        }
      }
    }, {
      '$sort': {
        'numScenes': -1
      }
    }
  ]
})

pandas.DataFrame(list(res))
```

Out[ ]:

	seasonNum	episodeNum	episodeTitle	numScenes
0	8	3	The Long Night	292
1	8	5	The Bells	220
2	2	9	Blackwater	133
3	8	4	The Last of the Starks	113
4	7	7	The Dragon and the Wolf	104

M2

- Write and execute a query that shows the first three episodes for each season
- Your aggregation should have the following attributes:
  - seasonNum
  - firstThreeEpisodes , which is an array that contains the titles of the first, second, and third episodes (in that order) of the season
- Order your output on seasonNum ascending
  - It's okay if the firstThreeEpisodes column is a bit truncated by the dataframe

```
In [ ]: res = episodes.aggregate([
  {
    '$match': {
      'episodeNum': {
        '$lte': 3
      }
    }
  }, {
    '$group': {
      '_id': '$seasonNum',
      'firstThreeEpisodes': {
        '$push': '$episodeTitle'
      }
    }
  }, {
    '$sort': {
      '_id': 1
    }
  }, {
    '$project': {
      '_id': 0,
      'firstThreeEpisodes': 1,
      'seasonNum': '$_id'
    }
  }
])
```

```
])

pandas.DataFrame(list(res))
```

Out[ ]:

	firstThreeEpisodes	seasonNum
0	[Winter Is Coming, The Kingsroad, Lord Snow]	1
1	[The North Remembers, The Night Lands, What Is...	2
2	[Valar Dohaeris, Dark Wings, Dark Words, Walk ...	3
3	[Two Swords, The Lion and the Rose, Breaker of...	4
4	[The Wars to Come, The House of Black and Whit...	5
5	[The Red Woman, Home, Oathbreaker]	6
6	[Dragonstone, Stormborn, The Queen's Justice]	7
7	[Winterfell, A Knight of the Seven Kingdoms, T...	8

### M3

- Write and execute a query that shows statistics about each season
- Your aggregation should have the following attributes:
  - `seasonNum`
  - `numEpisodes` , which is the number of episodes in the season
  - `startDate` , which is the earliest air date associated with an episode in the season
  - `endDate` , which is the latest air date associated with an episode in the season
  - `shortestEpisodeLength`
  - `longestEpisodeLength`
    - The length of an episode is the greatest `sceneEnd` value in the episode's `scenes` array
- Order your output on `seasonNum` ascending

```
In [ ]: res = episodes.aggregate([
    {
        '$group': {
            '_id': '$seasonNum',
            'numEpisodes': {
                '$sum': 1
            },
            'startDate': {
                '$min': '$episodeAirDate'
            },
            'endDate': {
                '$max': '$episodeAirDate'
            },
            'shortestEpisodeLength': {
                '$min': {
                    '$max': '$scenes.sceneEnd'
                }
            },
            'longestEpisodeLength': {
                '$max': {
                    '$max': '$scenes.sceneEnd'
                }
            }
        }
    }, {
        '$project': {
```

```
        '_id': 0,
        'seasonNum': '$_id',
        'numEpisodes': 1,
        'startDate': 1,
        'endDate': 1,
        'shortestEpisodeLength': 1,
        'longestEpisodeLength': 1
    }, {
    }, {
    '$sort': {
        'seasonNum': 1
    }
}
])

pandas.DataFrame(list(res))
```

Out [ ]:

	numEpisodes	startDate	endDate	shortestEpisodeLength	longestEpisodeLength	seasonNum
0	10	2011-04-17	2011-06-19	0:51:30	1:00:57	1
1	10	2012-04-01	2012-06-03	0:49:18	1:02:04	2
2	10	2013-03-31	2013-06-09	0:49:25	1:01:20	3
3	10	2014-04-06	2014-06-15	0:49:19	1:04:49	4
4	10	2015-03-29	2015-06-14	0:50:32	1:02:01	5
5	10	2016-04-24	2016-06-26	0:51:52	1:10:14	6
6	7	2017-07-16	2017-08-27	0:50:05	1:21:10	7
7	6	2019-04-14	2019-05-19	0:54:29	1:21:37	8

## M4

- Write and execute a query that shows sublocations and the scenes they appear in
- Your aggregation should have the following attributes:
  - subLocation
  - totalScenes , which is the number of scenes that are set in the sublocation
  - firstSeasonNum
  - firstEpisodeNum
    - (firstSeasonNum, firstEpisodeNum) identifies the first episode that the sublocation appears in
  - lastSeasonNum
  - lastEpisodeNum
    - (lastSeasonNum, lastEpisodeNum) identifies the last episode that the sublocation appears in
- Order your output on totalScenes descending, and only keep the sublocations with more than 50 scenes

```
In [ ]: res = episodes.aggregate([
    {
        '$unwind': {
            'path': '$scenes',
            'preserveNullAndEmptyArrays': False
        }
    }, {
        '$group': {
            '_id': '$scenes.subLocation',
            'totalScenes': {
                '$sum': 1
            }
        },
```

```
        'firstAppear': {
            '$min': {
                'seasonNum': '$seasonNum',
                'episodeNum': '$episodeNum'
            }
        },
        'lastAppear': {
            '$max': {
                'seasonNum': '$seasonNum',
                'episodeNum': '$episodeNum'
            }
        }
    }, {
        '$project': {
            'subLocation': '$_id',
            '_id': 0,
            'totalScenes': 1,
            'firstSeasonNum': '$firstAppear.seasonNum',
            'firstEpisodeNum': '$firstAppear.episodeNum',
            'lastSeasonNum': '$lastAppear.seasonNum',
            'lastEpisodeNum': '$lastAppear.episodeNum'
        }
    }, {
        '$match': {
            'totalScenes': {
                '$gt': 50
            }
        }
    }, {
        '$sort': {
            'totalScenes': -1
        }
    }
])

pandas.DataFrame(list(res))
```

Out[ ]:

	totalScenes	subLocation	firstSeasonNum	firstEpisodeNum	lastSeasonNum	lastEpisodeNum
0	1094	King's Landing	1	1	8	6
1	734	Winterfell	1	1	8	6
2	494	None	1	1	8	6
3	267	Castle Black	1	1	8	6
4	142	Dragonstone	2	1	8	5
5	77	The Haunted Forest	1	1	8	6
6	69	Outside Winterfell	1	1	8	4
7	66	Craster's Keep	2	1	4	5
8	60	The Wall	2	10	8	6
9	57	The Twins	1	9	7	1
10	56	Blackwater Rush	7	4	7	5
11	53	Blackwater Bay	2	8	8	6



# Neo4j

- The cell below creates nodes and relationships that model movies and the people involved in them

```
In [ ]: with open("movies.txt") as f:
        queries = str(f.read())

graph.execute_query("match (p:Person), (m:Movie) detach delete p, m")
graph.execute_query(queries)
print("Successfully inserted movie data")
```

Successfully inserted movie data

## N1

- Write and execute a cypher that shows actors and the number of movies they appear in
  - You should focus only on the `ACTED_IN` relationship, no other relationship
- Your output should have the following attributes:
  - `name` , which is the name of the actor
  - `num_movies`
- Order your output on `num_movies` descending, and only keep actors who have acted in 4 or more movies

```
In [ ]: res = graph.execute_query("""
        match (a: Person)-[:ACTED_IN]->(m: Movie)
        with a.name as name, count(distinct m) as num_movies
        where num_movies >= 4
        return name, num_movies
        order by num_movies desc
        """)

pandas.DataFrame([dict(r) for r in res.records])
```

Out [ ]:

	name	num_movies
0	Tom Hanks	12
1	Keanu Reeves	7
2	Hugo Weaving	5
3	Jack Nicholson	5
4	Meg Ryan	5
5	Cuba Gooding Jr.	4

## N2

- Write and execute a cypher that shows people and movies they either acted in or directed
- Your output should have the following attributes:
  - `name` , which is the name of the person
  - `directed_movies` , which is an array of titles of movies that the person directed
  - `acted_in_movies` , which is an array of titles of movies that the person acted in
- Order your output on `name` ascending, and only keep people that have directed at least one movie **and** acted in at least one movie (i.e., there should be no empty arrays. Arrays with one element are fine.)

```
In [ ]: res = graph.execute_query("""
        match (p: Person)-[:ACTED_IN]-(m1: Movie), (p: Person)-[:DIRECTED]->(m2: Movie)
```

```
with p.name as name, collect(DISTINCT m1.title) as acted_in_movies, collect(m2.title) as directed_movies
where size(acted_in_movies) > 0 and size(directed_movies) > 0
return name, acted_in_movies, directed_movies
order by name asc
"""
)

pandas.DataFrame([dict(r) for r in res.records])
```

Out [ ]:

	name	acted_in_movies	directed_movies
0	Clint Eastwood	[Unforgiven]	[Unforgiven]
1	Danny DeVito	[Hoffa, One Flew Over the Cuckoo's Nest]	[Hoffa, Hoffa]
2	James Marshall	[A Few Good Men]	[V for Vendetta, Ninja Assassin]
3	Tom Hanks	[You've Got Mail, Sleepless in Seattle, Joe Ve...	[That Thing You Do, That Thing You Do, That Th...
4	Werner Herzog	[What Dreams May Come]	[RescueDawn]

### N3

- Write and execute a cypher that shows people and movies they both acted in and directed
- Your output should have the following attributes:
  - `name` , which is the name of the person
  - `acted_in_and_directed_movies` , which is an array of titles of movies that the person both acted in and directed
- Order your output on `name` ascending, and only keep people that have acted in at least one movie that they directed (i.e., there should be no empty arrays. Arrays with one element are fine.)

```
In [ ]: res = graph.execute_query("""
match (p: Person)-[:ACTED_IN]-(m1: Movie), (p: Person)-[:DIRECTED]->(m2: Movie)
with p.name as name, collect(DISTINCT m1.title) as acted_in_movies, collect(m2.title) as directed_movies
where size([i in acted_in_movies where i in directed_movies]) > 0
return name, [i in acted_in_movies where i in directed_movies] as acted_in_and_directed_movies
order by name asc
""")

pandas.DataFrame([dict(r) for r in res.records])
```

Out [ ]:

	name	acted_in_and_directed_movies
0	Clint Eastwood	[Unforgiven]
1	Danny DeVito	[Hoffa]
2	Tom Hanks	[That Thing You Do]

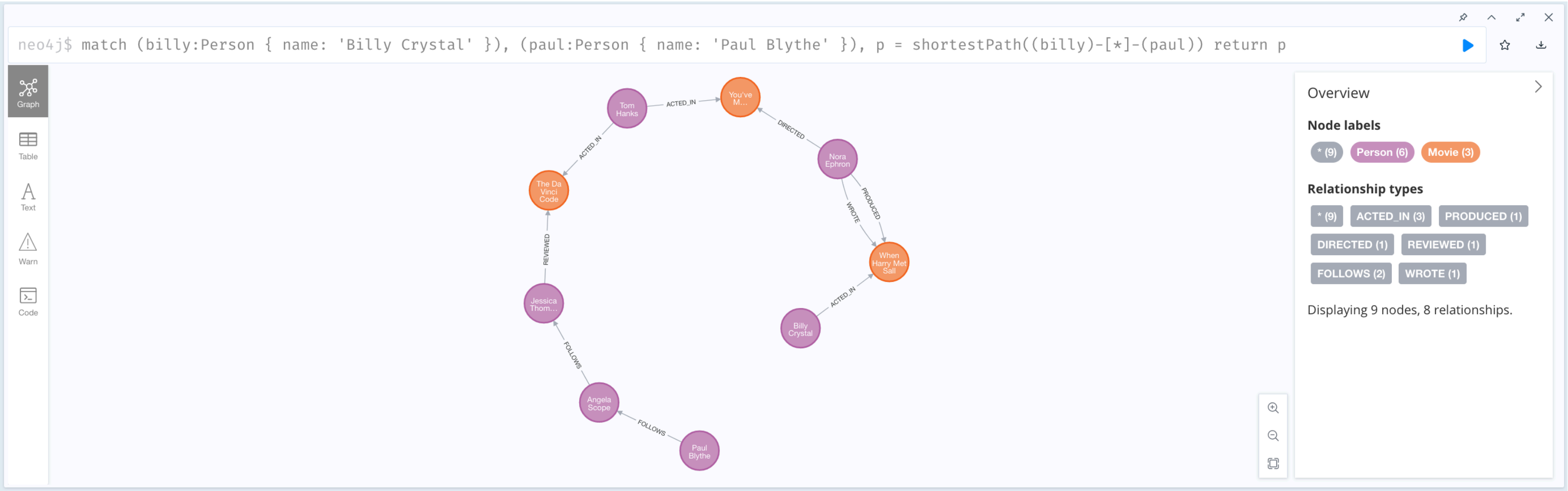
### N4

- Write and execute a cypher that shows pairs of people and how closely connected they are
- Your output should have the following attributes:
  - `person_1_name` , which is the name of the first person in the pair
  - `person_2_name` , which is the name of the second person in the pair
  - `num_people_between` , which is the number of people (including the pair itself) separating the pair. You should use the `shortestPath` function to compute this.
- To prevent duplicates in your output, you should only keep rows where `person_1_name < person_2_name`
- Order your output on `(person_1_name, person_2_name)` , and only keep rows where `num_people_between > 5`

- As an example, you should get the following row in your output:

person_1_name	person_2_name	num_people_between
Billy Crystal	Paul Blythe	6

- The shortest path between Billy Crystal and Paul Blythe is shown below
  - `num_people_between` is 6 because there are 6 nodes marked as `Person` (including Billy's and Paul's nodes)



```
In [ ]: res = graph.execute_query("""
match (p1: Person), (p2: Person)
where p1.name < p2.name
match path = shortestPath((p1)-[*]-(p2))
with p1.name as person_1_name, p2.name as person_2_name, size([i in nodes(path) where i:Person]) as num_people_between
where num_people_between > 5
return person_1_name, person_2_name, num_people_between
order by person_1_name, person_2_name
""")

pandas.DataFrame([dict(r) for r in res.records])
```

Out [ ]:

	person_1_name	person_2_name	num_people_between
0	Billy Crystal	Paul Blythe	6
1	Bruno Kirby	Paul Blythe	6
2	Carrie Fisher	Paul Blythe	6
3	Christian Bale	Dina Meyer	6
4	Christian Bale	Ice-T	6
5	Christian Bale	Paul Blythe	6
6	Christian Bale	Robert Longo	6
7	Christian Bale	Takeshi Kitano	6
8	Ethan Hawke	Paul Blythe	6
9	Jan de Bont	Paul Blythe	6
10	Paul Blythe	Scott Hicks	6
11	Paul Blythe	Zach Grenier	6

## SQL To NoSQL

- You will move relational data to document and graph databases
  - You will do your modeling in Python. You shouldn't be writing any SQL.**
- You will be using the `classicmodels` database for this section. You may want to drop the database and re-run the SQL script (included in the directory) to ensure you have the right data.
  - You will be modeling customers and the products they ordered

## MongoDB: Customers

- For the document database, you will create two collections: `customers` and `products`
- `customers` will contain customer information as well as all the orders they've placed
- You will use `customer_orders_all_df` to create your `customers` collection

In [ ]:

```
%%sql

customer_orders_all <<

select
    c.customerNumber, c.customerName, c.country, o.orderNumber,
    o.orderDate, od.productCode, od.quantityOrdered, od.priceEach
from classicmodels.orders o
join classicmodels.customers c using (customerNumber)
join classicmodels.orderdetails od using (orderNumber);

* mysql+pymysql://root:***@localhost
2996 rows affected.
Returning data to local variable customer_orders_all
```

In [ ]:

```
customer_orders_all_df = customer_orders_all.DataFrame()
customer_orders_all_df.head(10)
```

Out[ ]:

	customerNumber	customerName	country	orderNumber	orderDate	productCode	quantityOrdered	priceEach
0	363	Online Diecast Creations Co.	USA	10100	2003-01-06	S18_1749	30	136.00
1	363	Online Diecast Creations Co.	USA	10100	2003-01-06	S18_2248	50	55.09
2	363	Online Diecast Creations Co.	USA	10100	2003-01-06	S18_4409	22	75.46
3	363	Online Diecast Creations Co.	USA	10100	2003-01-06	S24_3969	49	35.29
4	128	Blauer See Auto, Co.	Germany	10101	2003-01-09	S18_2325	25	108.06
5	128	Blauer See Auto, Co.	Germany	10101	2003-01-09	S18_2795	26	167.06
6	128	Blauer See Auto, Co.	Germany	10101	2003-01-09	S24_1937	45	32.53
7	128	Blauer See Auto, Co.	Germany	10101	2003-01-09	S24_2022	46	44.35
8	181	Vitachrome Inc.	USA	10102	2003-01-10	S18_1342	39	95.55
9	181	Vitachrome Inc.	USA	10102	2003-01-10	S18_1367	41	43.13

- Below is an example of how a customer and their orders are stored in MySQL, and how the document should look like in MongoDB
- The document should have the following attributes:
  - customerNumber
  - customerName
  - country
  - orders , which is a list of objects. Each object represents one order
    - orderNumber
    - orderDate
    - orderContents , which is a list of objects. Each object represents one product in the order
      - productCode
      - quantityOrdered
      - priceEach

MySQL relation:

	customerNumber	customerName	country	orderNumber	orderDate	productCode	quantityOrdered	priceEach
0	103	Atelier graphique	France	10123	2003-05-20	S18_1589	26	120.71
1	103	Atelier graphique	France	10123	2003-05-20	S18_2870	46	114.84
2	103	Atelier graphique	France	10123	2003-05-20	S18_3685	34	117.26
3	103	Atelier graphique	France	10123	2003-05-20	S24_1628	50	43.27
4	103	Atelier graphique	France	10298	2004-09-27	S10_2016	39	105.86
5	103	Atelier graphique	France	10298	2004-09-27	S18_2625	32	60.57
6	103	Atelier graphique	France	10345	2004-11-25	S24_2022	43	38.98

MongoDB document:

```
{
  customerNumber: 103
  customerName: "Atelier graphique",
  country: "France",
  orders: [
    {
      orderNumber: 10123,
```

```

        orderDate: "2003-05-20",
        orderContents: [
            {
                productCode: "S18_1589",
                quantityOrdered: 26,
                priceEach: "120.71"
            },
            {
                productCode: "S18_2870",
                quantityOrdered: 46,
                priceEach: "114.84"
            },
            {
                productCode: "S18_3685",
                quantityOrdered: 34,
                priceEach: "117.26"
            },
            {
                productCode: "S24_1628",
                quantityOrdered: 50,
                priceEach: "43.27"
            }
        ]
    },
    {
        orderNumber: 10298,
        orderDate: "2004-09-27",
        orderContents: [
            {
                productCode: "S10_2016",
                quantityOrdered: 39,
                priceEach: "105.86"
            },
            {
                productCode: "S18_2625",
                quantityOrdered: 32,
                priceEach: "60.57"
            }
        ]
    },
    {
        orderNumber: 10345,
        orderDate: "2004-11-25",
        orderContents: [
            {
                productCode: "S24_2022",
                quantityOrdered: 43,
                priceEach: "38.98"
            }
        ]
    }
]
}

```

In [ ]: *# TODO: Create a list of dicts. Each dict represents one customer.*

"""

Tips:

To iterate through dataframe:

```

for _, r in customer_orders_all_df.iterrows():
    r = dict(r)
    Access fields like r['customerName'], r['country'], ...

```

The orderDate and priceEach fields are stored as datetime.date and Decimal objects in the dataframe. These types are not compatible with the pymongo API. You can convert them to strings by calling str(r['orderDate']) and str(r['priceEach']). Alternatively, you can look into the datetime.datetime and bson.decimal128.Decimal128 objects, which are supported by pymongo.

```

"""

```

```

dic = dict()

```

```

for _, r in customer_orders_all_df.iterrows():
    r = dict(r)
    # print(r)

    # customer info
    if r['customerNumber'] not in dic.keys():
        cus_info = {}
        cus_info['customerNumber'] = r['customerNumber']
        cus_info['customerName'] = r['customerName']
        cus_info['country'] = r['country']
        cus_info['orders'] = {}
        dic[r['customerNumber']] = cus_info

    # orders info
    if r['orderNumber'] not in dic[r['customerNumber']]['orders'].keys():
        order_info = {}
        order_info['orderNumber'] = r['orderNumber']
        order_info['orderDate'] = str(r['orderDate'])
        order_info['orderContents'] = []
        dic[r['customerNumber']]['orders'][r['orderNumber']] = order_info

    # product info - we always need product info
    pro_info = {}
    pro_info['productCode'] = r['productCode']
    pro_info['quantityOrdered'] = r['quantityOrdered']
    pro_info['priceEach'] = str(r['priceEach'])
    dic[r['customerNumber']]['orders'][r['orderNumber']]['orderContents'].append(pro_info)

# Delete the order number keys
for i in dic.keys():
    dic[i]['orders'] = list(dic[i]['orders'].values())

# Delete the customer number keys
your_list_of_customers = list(dic.values())

# print(dic)

```

```

In [ ]: def insert_customers(d):
        mongo_client['w4111']['customers'].drop()
        mongo_client['w4111']['customers'].insert_many(d)

# TODO: Put the name of your list of dicts below
insert_customers(your_list_of_customers)
print("Successfully inserted customer data")

```

Successfully inserted customer data

# MongoDB: Products

- To create the `products` collection, you will use `products_all_df`
- A document in `products` simply contains product information, as shown below

```
{
  productCode: "S10_1678",
  productName: "1969 Harley Davidson Ultimate Chopper",
  productVendor: "Min Lin Diecast"
}
```

```
In [ ]: %%sql

products_all <<

select productCode, productName, productVendor
from classicmodels.products;

* mysql+pymysql://root:***@localhost
110 rows affected.
Returning data to local variable products_all
```

```
In [ ]: products_all_df = products_all.DataFrame()
products_all_df.head(10)
```

Out [ ]:

	productCode	productName	productVendor
0	S10_1678	1969 Harley Davidson Ultimate Chopper	Min Lin Diecast
1	S10_1949	1952 Alpine Renault 1300	Classic Metal Creations
2	S10_2016	1996 Moto Guzzi 1100i	Highway 66 Mini Classics
3	S10_4698	2003 Harley-Davidson Eagle Drag Bike	Red Start Diecast
4	S10_4757	1972 Alfa Romeo GTA	Motor City Art Classics
5	S10_4962	1962 LanciaA Delta 16V	Second Gear Diecast
6	S12_1099	1968 Ford Mustang	Autoart Studio Design
7	S12_1108	2001 Ferrari Enzo	Second Gear Diecast
8	S12_1666	1958 Setra Bus	Welly Diecast Productions
9	S12_2823	2002 Suzuki XREO	Unimax Art Galleries

```
In [ ]: # TODO: Create a list of dicts. Each dict represents one product.

"""
Tips:

    To iterate through dataframe:

        for _, r in products_all_df.iterrows():
            r = dict(r)
            Access fields like r['productName'], r['productVendor'], ...

"""

your_list_of_products = []

for _, r in products_all_df.iterrows():
    r = dict(r)
```



```
your_list_of_products.append(r)
```

```
In [ ]: def insert_products(d):
        mongo_client['w4111']['products'].drop()
        mongo_client['w4111']['products'].insert_many(d)

# TODO: Put the name of your list of dicts below
insert_products(your_list_of_products)
print("Successfully inserted product data")
```

Successfully inserted product data

## MongoDB: Testing

- Run through the following cells
- **Make sure the outputs are completely visible. You shouldn't need to scroll to see the entire output.**
  - You may need to click on the blank section immediately to the left of your output to toggle between scrolling and unscrolling

```
In [ ]: import json

def prepr(doc):
    try:
        del doc['_id']
    except KeyError:
        pass

    def convert_str(d):
        if isinstance(d, dict):
            for k, v in d.items():
                d[k] = convert_str(v)
            return d
        elif isinstance(d, list):
            for i, v in enumerate(d):
                d[i] = convert_str(v)
            return d
        else:
            return str(d)

    convert_str(doc)
    return json.dumps(doc, indent=2)
```

```
In [ ]: res = mongo_client['w4111']['customers'].aggregate([
        {
            '$match': {
                'customerNumber': 219
            }
        }
    ])

print(prepr(list(res)[0]))
```

```
{
  "customerNumber": "219",
  "customerName": "Boards & Toys Co.",
  "country": "USA",
  "orders": [
    {
      "orderNumber": "10154",
      "orderDate": "2003-10-02",
      "orderContents": [
        {
          "productCode": "S24_3151",
          "quantityOrdered": "31",
          "priceEach": "75.23"
        },
        {
          "productCode": "S700_2610",
          "quantityOrdered": "36",
          "priceEach": "59.27"
        }
      ]
    },
    {
      "orderNumber": "10376",
      "orderDate": "2005-02-08",
      "orderContents": [
        {
          "productCode": "S12_3380",
          "quantityOrdered": "35",
          "priceEach": "98.65"
        }
      ]
    }
  ]
}
```

```
In [ ]: res = mongo_client['w4111']['customers'].aggregate([
    {
        '$match': {
            'customerNumber': 103
        }
    }
])

print(prepr(list(res)[0]))
```

```

{
  "customerNumber": "103",
  "customerName": "Atelier graphique",
  "country": "France",
  "orders": [
    {
      "orderNumber": "10123",
      "orderDate": "2003-05-20",
      "orderContents": [
        {
          "productCode": "S18_1589",
          "quantityOrdered": "26",
          "priceEach": "120.71"
        },
        {
          "productCode": "S18_2870",
          "quantityOrdered": "46",
          "priceEach": "114.84"
        },
        {
          "productCode": "S18_3685",
          "quantityOrdered": "34",
          "priceEach": "117.26"
        },
        {
          "productCode": "S24_1628",
          "quantityOrdered": "50",
          "priceEach": "43.27"
        }
      ]
    },
    {
      "orderNumber": "10298",
      "orderDate": "2004-09-27",
      "orderContents": [
        {
          "productCode": "S10_2016",
          "quantityOrdered": "39",
          "priceEach": "105.86"
        },
        {
          "productCode": "S18_2625",
          "quantityOrdered": "32",
          "priceEach": "60.57"
        }
      ]
    },
    {
      "orderNumber": "10345",
      "orderDate": "2004-11-25",
      "orderContents": [
        {
          "productCode": "S24_2022",
          "quantityOrdered": "43",
          "priceEach": "38.98"
        }
      ]
    }
  ]
}

```

```

In [ ]: res = mongo_client['w4111']['products'].aggregate([
  {
    '$match': {

```

```
        'productCode': 'S18_1889'
    }
}

print(prepr(list(res)[0]))
```

```
{
  "productCode": "S18_1889",
  "productName": "1948 Porsche 356-A Roadster",
  "productVendor": "Gearbox Collectibles"
}
```

## Neo4j: All Data

- For the graph database, you will have two types of nodes: `Customer` and `Product`
  - **Make sure to use these exact names**
- An order is represented as a relationship from the `Customer` node to the `Product` node. The type of the relationship should be `ORDERED`.
- You will use `customer_orders_limit_df` to create your graph

```
In [ ]: %%sql

customer_orders_limit <<

select
  c.customerNumber, c.customerName, c.country, o.orderNumber,
  o.orderDate, od.productCode, p.productName, p.productVendor,
  od.quantityOrdered, od.priceEach
from classicmodels.orders o
  join classicmodels.customers c using (customerNumber)
  join classicmodels.orderdetails od using (orderNumber)
  join classicmodels.products p using (productCode)
where od.quantityOrdered > 49;

* mysql+pymysql://root:***@localhost
139 rows affected.
Returning data to local variable customer_orders_limit
```

```
In [ ]: customer_orders_limit_df = customer_orders_limit.DataFrame()
customer_orders_limit_df.head(10)
```

Out[ ]:

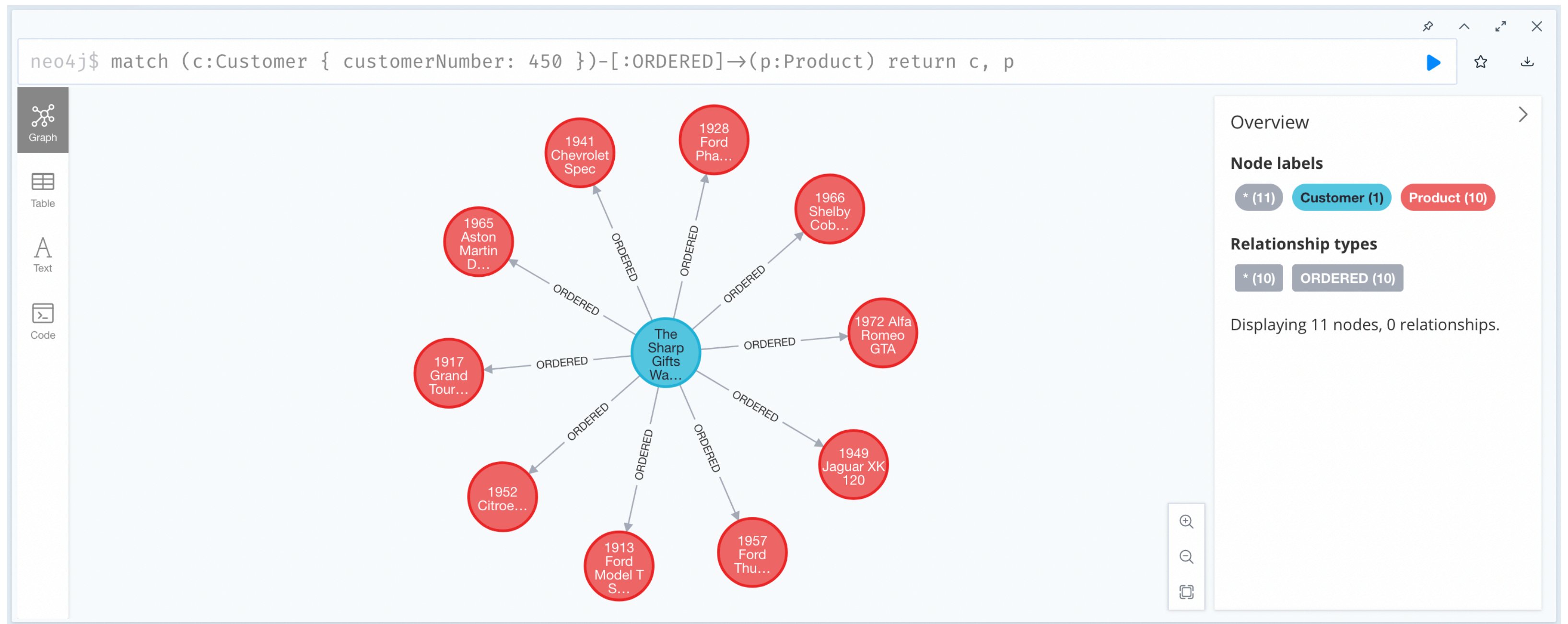
	customerNumber	customerName	country	orderNumber	orderDate	productCode	productName	productVendor	quantityOrdered	priceEach
0	363	Online Diecast Creations Co.	USA	10100	2003-01-06	S18_2248	1911 Ford Town Car	Motor City Art Classics	50	55.09
1	145	Danish Wholesale Imports	Denmark	10105	2003-02-11	S10_4757	1972 Alfa Romeo GTA	Motor City Art Classics	50	127.84
2	145	Danish Wholesale Imports	Denmark	10105	2003-02-11	S24_3816	1940 Ford Delivery Sedan	Carousel DieCast Legends	50	75.47
3	278	Rovelli Gifts	Italy	10106	2003-02-17	S24_3949	Corsair F4U ( Bird Cage)	Second Gear Diecast	50	55.96
4	124	Mini Gifts Distributors Ltd.	USA	10113	2003-03-26	S18_4668	1939 Cadillac Limousine	Studio M Art Models	50	43.27
5	148	Dragon Souveniers, Ltd.	Singapore	10117	2003-04-16	S72_3212	Pont Yacht	Unimax Art Galleries	50	52.42
6	353	Reims Collectables	France	10121	2003-05-07	S12_2823	2002 Suzuki XREO	Unimax Art Galleries	50	126.52
7	103	Atelier graphique	France	10123	2003-05-20	S24_1628	1966 Shelby Cobra 427 S/C	Carousel DieCast Legends	50	43.27
8	458	Corrida Auto Replicas, Ltd	Spain	10126	2003-05-28	S18_4600	1940s Ford truck	Motor City Art Classics	50	102.92
9	324	Stylish Desk Decors, Co.	UK	10129	2003-06-12	S24_3816	1940 Ford Delivery Sedan	Carousel DieCast Legends	50	76.31

- Below is an example of how a customer and their orders are stored in MySQL, and how the graph should look like in Neo4j
  - Note that the same order may be represented as many relationships since one order could contain many products
- The `Customer` nodes should have the following attributes:
  - `customerNumber`
  - `customerName`
  - `country`
- The `Product` nodes should have the following attributes:
  - `productCode`
  - `productName`
  - `productVendor`
- The `ORDERED` relationships should have the following attributes:
  - `orderNumber`
  - `orderDate`
  - `quantityOrdered`
  - `priceEach`

MySQL relation:

	customerNumber	customerName	country	orderNumber	orderDate	productCode	productName	productVendor	quantityOrdered	priceEach
0	450	The Sharp Gifts Warehouse	USA	10250	2004-05-11	S32_4289	1928 Ford Phaeton Deluxe	Highway 66 Mini Classics	50	62.6
1	450	The Sharp Gifts Warehouse	USA	10257	2004-06-14	S18_2949	1913 Ford Model T Speedster	Carousel DieCast Legends	50	92.19
2	450	The Sharp Gifts Warehouse	USA	10400	2005-04-01	S10_4757	1972 Alfa Romeo GTA	Motor City Art Classics	64	134.64
3	450	The Sharp Gifts Warehouse	USA	10400	2005-04-01	S18_3856	1941 Chevrolet Special Deluxe Cabriolet	Exoto Designs	58	88.93
4	450	The Sharp Gifts Warehouse	USA	10407	2005-04-22	S18_1589	1965 Aston Martin DB5	Classic Metal Creations	59	114.48
5	450	The Sharp Gifts Warehouse	USA	10407	2005-04-22	S18_1749	1917 Grand Touring Sedan	Welly Diecast Productions	76	141.1
6	450	The Sharp Gifts Warehouse	USA	10407	2005-04-22	S18_4933	1957 Ford Thunderbird	Studio M Art Models	66	64.14
7	450	The Sharp Gifts Warehouse	USA	10407	2005-04-22	S24_1628	1966 Shelby Cobra 427 S/C	Carousel DieCast Legends	64	45.78
8	450	The Sharp Gifts Warehouse	USA	10407	2005-04-22	S24_2766	1949 Jaguar XK 120	Classic Metal Creations	76	81.78
9	450	The Sharp Gifts Warehouse	USA	10407	2005-04-22	S24_2887	1952 Citroen-15CV	Exoto Designs	59	98.65

Neo4j graph:



```
In [ ]: # Deletes all customers and products (and their relationships).
# Feel free to run this as many times as you want to reset your data.
_ = graph.execute_query("""
    match (c:Customer), (p:Product)
    detach delete c, p
    """)
```

```
In [ ]: # TODO: Write and execute queries to create nodes and relationships
```

"""

Tips:

To iterate through dataframe:

```
for _, r in customer_orders_limit_df.iterrows():
    r = dict(r)
    Access fields like r['customerName'], r['country'], ...
```

The priceEach field are stored as a Decimal object in the dataframe. This type is not compatible with the neo4j API. You can convert it to a string by calling `str(r['priceEach'])`.

You should call `graph.execute_query` to execute your queries. This method takes in a second optional argument, a dict. This allows you to do query parameters. For instance, to execute the query in the screenshot above, you could run

```

graph.execute_query(
    "match (c:Customer { customerNumber: $custNum })-[:ORDERED]->(p:Product) return c, p",
    { "custNum": 450 }
)

"""

customer = {}
product = {}
order = []

for _, r in customer_orders_limit_df.iterrows():
    r = dict(r)

    # customer
    if r['customerNumber'] not in customer.keys():
        c = {}
        c['customerNumber'] = r['customerNumber']
        c['customerName'] = r['customerName']
        c['country'] = r['country']
        customer[r['customerNumber']] = c

    # product
    if r['productCode'] not in product.keys():
        p = {}
        p['productCode'] = r['productCode']
        p['productName'] = r['productName']
        p['productVendor'] = r['productVendor']
        product[r['productCode']] = p

    # order
    o = {}
    o['orderNumber'] = r['orderNumber']
    o['orderDate'] = str(r['orderDate'])
    o['quantityOrdered'] = r['quantityOrdered']
    o['priceEach'] = str(r['priceEach'])
    o['customerNumber'] = r['customerNumber']
    o['productCode'] = r['productCode']
    order.append(o)

# Delete the keys of customer and order
customer = list(customer.values())
product = list(product.values())

```

In [ ]: *# Declare: following codes are written by discussion with GPT, because I don't know how to insert data into neo4j*

```

def create_customer(tx, customer):
    query = (
        "CREATE (c:Customer {customerNumber: $customerNumber, customerName: $customerName, country: $country})"
    )
    tx.run(query, customer)

def create_product(tx, product):
    query = (
        "CREATE (p:Product {productCode: $productCode, productName: $productName, productVendor: $productVendor})"
    )
    tx.run(query, product)

def create_order(tx, order):
    query = (
        "MATCH (c:Customer {customerNumber: $customerNumber}), "
        "      (p:Product {productCode: $productCode}) "
        "CREATE (c)-[o:ORDERED {orderNumber: $orderNumber, orderDate: $orderDate, "
        "quantityOrdered: $quantityOrdered, priceEach: $priceEach}]->(p)"
    )

```



```
)
tx.run(query, order)

# create nodes
with graph.session() as session:
    for c in customer:
        session.write_transaction(create_customer, c)
    for p in product:
        session.write_transaction(create_product, p)

# create edges
with graph.session() as session:
    for o in order:
        session.write_transaction(create_order, o)
```

/var/folders/4m/rxd9grt9lhq7n4\_\_bqp6nkvr0000gn/T/ipykernel\_36129/51608906.py:27: DeprecationWarning: write\_transaction has been renamed to execute\_write  
session.write\_transaction(create\_customer, c)  
/var/folders/4m/rxd9grt9lhq7n4\_\_bqp6nkvr0000gn/T/ipykernel\_36129/51608906.py:29: DeprecationWarning: write\_transaction has been renamed to execute\_write  
session.write\_transaction(create\_product, p)  
/var/folders/4m/rxd9grt9lhq7n4\_\_bqp6nkvr0000gn/T/ipykernel\_36129/51608906.py:34: DeprecationWarning: write\_transaction has been renamed to execute\_write  
session.write\_transaction(create\_order, o)

## Neo4j: Testing

- Run through the following cells
- **Make sure the outputs are fully visible**

In [ ]:

```
res = []

for r in graph.execute_query("""
    match (c:Customer { customerNumber: 412 })-[o:ORDERED]->(p:Product)
    return c, o, p
""").records:
    res.append(dict(r['c'])) | dict(r['o']) | dict(r['p']))

pandas.DataFrame(res)
```

Out [ ]:

	country	customerNumber	customerName	orderNumber	quantityOrdered	orderDate	priceEach	productCode	productName	productVendor
0	New Zealand	412	Extreme Desk Decorations, Ltd	10418	52	2005-05-16	64.41	S24_2360	1982 Ducati 900 Monster	Highway 66 Mini Classics
1	New Zealand	412	Extreme Desk Decorations, Ltd	10234	50	2004-03-30	146.65	S18_1662	1980s Black Hawk Helicopter	Red Start Diecast
2	New Zealand	412	Extreme Desk Decorations, Ltd	10268	50	2004-07-12	124.59	S18_2325	1932 Model A Ford J-Coupe	Autoart Studio Design
3	New Zealand	412	Extreme Desk Decorations, Ltd	10418	50	2005-05-16	100.01	S32_4485	1974 Ducati 350 Mk3 Desmo	Second Gear Diecast

In [ ]:

```
res = []

for r in graph.execute_query("""
    match (c:Customer)-[o:ORDERED]->(p:Product { productCode: 'S12_2823' })
    return c, o, p
""").records:
    res.append(dict(r['c'])) | dict(r['o']) | dict(r['p']))

pandas.DataFrame(res)
```



Out[ ]:

	country	customerNumber	customerName	orderNumber	quantityOrdered	orderDate	priceEach	productCode	productName	productVendor
0	France	353	Reims Collectables	10121	50	2003-05-07	126.52	S12_2823	2002 Suzuki XREO	Unimax Art Galleries
1	Austria	382	Salzburg Collectables	10341	55	2004-11-24	120.50	S12_2823	2002 Suzuki XREO	Unimax Art Galleries
2	UK	201	UK Collectables, Ltd.	10403	66	2005-04-08	122.00	S12_2823	2002 Suzuki XREO	Unimax Art Galleries

In [ ]:

```
res = []

for r in graph.execute_query("""
    match (c:Customer)-[o:ORDERED { quantityOrdered: 60 }]->(p:Product)
    return c, o, p
""").records:
    res.append(dict(r['c']) | dict(r['o']) | dict(r['p']))

pandas.DataFrame(res)
```

Out[ ]:

	country	customerNumber	customerName	orderNumber	quantityOrdered	orderDate	priceEach	productCode	productName	productVendor
0	Spain	141	Euro+ Shopping Channel	10412	60	2005-05-03	157.49	S18_3232	1992 Ferrari 360 Spider red	Unimax Art Galleries
1	USA	362	Gifts4AllAges.com	10414	60	2005-05-06	72.58	S24_3151	1912 Ford Model T Delivery Wagon	Min Lin Diecast
2	Australia	282	Souvenirs And Things Co.	10420	60	2005-05-29	60.26	S24_1046	1970 Chevy Chevelle SS 454	Unimax Art Galleries

In [ ]:

```
res = []

for r in graph.execute_query("""
    match (n:Customer | Product)
    return labels(n) as type, count(*) as count
union
    match ()-[r:ORDERED]->()
    return type(r) as type, count(*) as count
""").records:
    res.append(dict(r))

pandas.DataFrame(res)
```

Out[ ]:

	type	count
0	[Customer]	57
1	[Product]	84
2	ORDERED	139