

BASES DE DATOS NO ESTRUCTURADAS

UNAM— Semestre 2022-2

Integrantes:
Andrés Urbano Guillermo Gerardo
Avitúa Varela Fernando
Santa Rita Vizuet Fernando

§ Bases de Datos orientadas a grafos §

Introducción

Neo4j es un software libre de Base de datos orientada a grafos, implementado en Java. Los desarrolladores describen a *Neo4j* como un motor de persistencia embebido, basado en disco, implementado en Java, completamente transaccional, que almacena datos estructurados en grafos en lugar de en tablas ("embedded, disk-based, fully transactional Java persistence engine that stores data structured in graphs rather than in tables"). La versión 1.0 de Neo4j fue lanzada en febrero de 2010.

En *Neo4j*, todo se almacena en forma de arista, nodo o atributo. Cada nodo y arista puede tener cualquier número de atributos y tanto los nodos como aristas se pueden etiquetar. Las etiquetas se pueden utilizar para limitar las búsquedas. A partir de la versión 2.0, se agregó la indexación a Cypher con la introducción de esquemas. Anteriormente, los índices se admitían por separado de Cypher.

Diseño e implementación de la base de datos

Una de las virtudes de las bases de datos orientadas a grafos es que pueden ser diseñadas vía dibujo, donde podemos ver las entidades y las relaciones que mantienen estas. El siguiente esquema muestra, de forma general, la idea detrás del diseño de la base de datos.

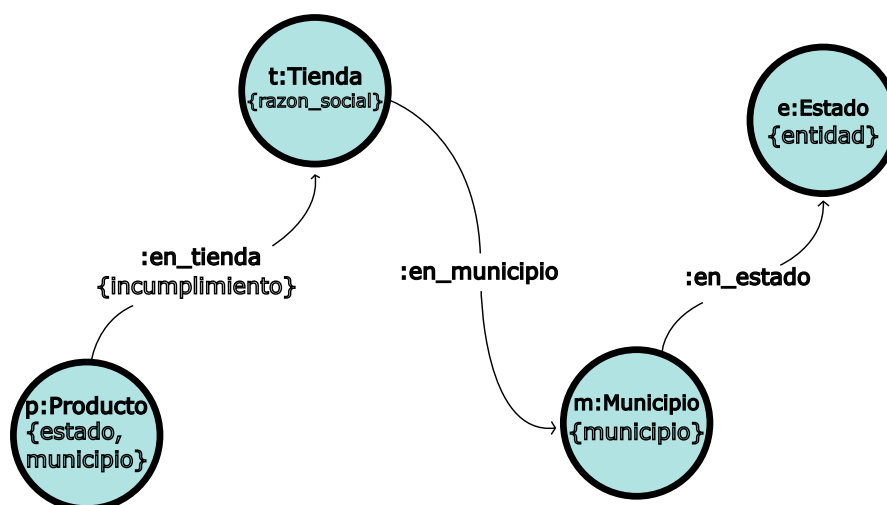


Figure 1: Diseño general de la base de datos.

El diseño de la gráfica de la figura 1 nos dice de forma natural la pertenencia entre cada una de las entidades. Existen atributos en el nodo de *Producto* relativos a la pertenencia del producto en cada uno de los estados y municipios del país; el agregar estos atributos ayuda a realizar consultas más sencillas, ayudando a satisfacer de mejor manera las necesidades del usuario final.

Dado el conjunto de datos a trabajar (**PROFECO**), cargamos este al un DataFrame de la librería Pandas y luego es exportado a la base de datos de *Neo4j*. Para ello, primero establecemos conexión con *Neo4j*:

```

1  # conexión con la base de datos
2  data_base_connection = GraphDatabase.driver(uri="bolt://localhost:7687",
3                                              auth=("neo4j", "1234"))
4  session = data_base_connection.session()

```

Con una conexión exitosa, es posible ejecutar las instrucciones para cargar el conjunto de datos a Neo4j:

```

1  # valores del DataFrame como listas de listas
2  profeco_lista = profeco_df.values.tolist()
3
4  # construcción de lista de instrucciones para Neo4j
5  instrucciones_neo4j = []
6
7  for transaccion in profeco_lista[:]:
8      instruccion=f'''
9  MERGE (e:Estado {{nombre_estado:"{transaccion[10]}"}})
10 MERGE (tp:TipoProcuto {{tipo: "{transaccion[5]}"}})
11 MERGE (t:Tienda {{razon_social:"{transaccion[7]}"}})
12 MERGE (m:Municipio {{nombre_municipio:"{transaccion[9]}"}})
13 MERGE (p:Producto {{nombre_producto:"{transaccion[6]}", estado:e.nombre_estado,
14                      municipio:m.nombre_municipio}})
15 MERGE (p)-[et:en_tienda {{incumplimiento:"{transaccion[11]}"}}]->(t)
16 MERGE (t)-[:en_municipio]->(m)
17 MERGE (m)-[:en_estado]->(e)
18 MERGE (p)-[:es_de_tipo]->(tp)
19 '''.replace('\n', ' ')
20
21      instrucciones_neo4j.append(instruccion)
22
23  # ejecutar instrucciones de instrucciones_neo4j
24  for i in tqdm(instrucciones_neo4j):
25      session.run(i)

```

Consultas

Consulta 1: Dado un estado y un producto, buscar lugares donde pueda encontrarlo.

Para realizar esta consulta basta encontrar todas las tiendas que están asociadas a un producto, cuidando que el estado y municipio estén conectados

```

1  def estado_producto(estado, producto):
2      '''
3      Regresa la tienda y municipio que tiene ese producto
4      '''
5      query=f'''
6  MATCH(m:Municipio)

```

```

7 MATCH(e:Estado {{nombre_estado:"{estado}"}})
8 MATCH(t:Tienda)
9 MATCH(p:Producto {{nombre_producto:"{producto}"
10 ,estado:"{estado}",municipio:m.nombre_municipio}})
11 MATCH (p)-[:en_tienda]->(t)
12 MATCH (m)-[:en_estado]-> (e)
13 MATCH (t)-[:en_municipio]->(m)
14 RETURN t,m'''
15 q_res=session.run(query)
16 return q_res.data()

```

Por ejemplo, haciendo el llamado de la función con los siguientes valores:

```

1 estado='baja california'
2 producto='productos con informacion comercial'

```

obtenemos

```

Producto: productos con informacion comercial
-----
Tienda:mercado y carniceria la fiesta
Municipio:tijuana
-----
Tienda:auto servicio ojeda
Municipio:tijuana
-----
Tienda:fruteria y dulceria huicho
Municipio:tijuana
-----
Tienda:dulceria arcoiris
Municipio:tijuana

```

Figure 2: Resultado de la consulta 1 para un producto

También lo podemos ver directamente en la interfaz gráfica de Neo4j.

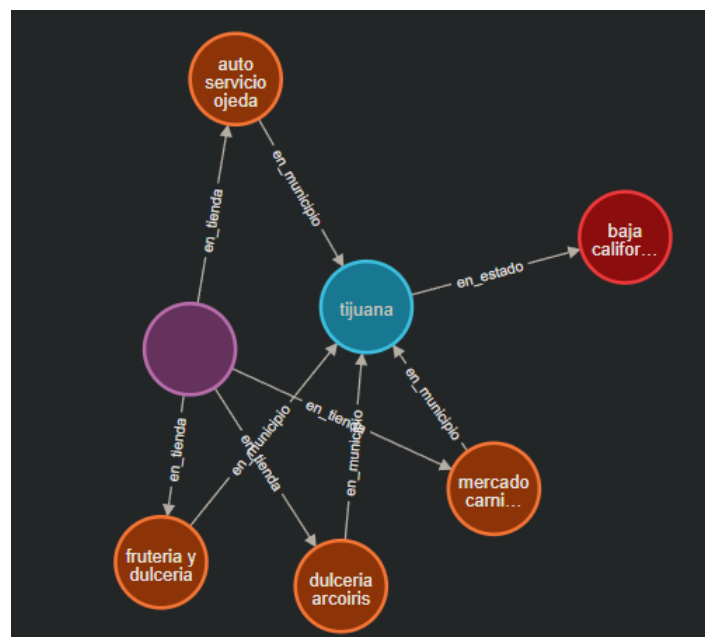


Figure 3: Resultado de la consulta 1.

Consulta 2: Dado un estado y una tienda, verificar si tiene algún incumplimiento con un producto.

El programa que resuelve esta consulta recibe como parámetros la tienda y estado. Si la salida es vacía, significa que en la tienda de dicho estado no hay ninguna clase de incumplimiento; de otro modo, serán listados los artículos que presentan dicho incumplimiento.

```

1  def falta_producto(tienda, estado):
2
3      query = f"""
4      MATCH (m:Municipio)
5      MATCH (p:Producto {{estado:"{estado}", municipio:m.nombre_municipio}})
6      MATCH (e:Estado {{nombre_estado:"{estado}"}})
7      MATCH (t:Tienda {{razon_social:"{tienda}"}})
8      MATCH (p) -[et:en_tienda]-> (t)
9      WHERE NOT et.incumplimiento = "no se detecto incumplimiento"
10     MATCH (t) -[:en_municipio]-> (m)
11     MATCH (m) -[:en_estado]-> (e)
12     RETURN p
13     """.replace('\n', ' ')
14
15     return session.run(query).data()

```

Obsérvese que se aprovecha la existencia de los atributos en los nodos, lo cuál permite realizar consultas sencillas. Como ejemplo, veremos los incumplimiento que se encuentran en el *Hotel del Fresno de Fresnillo S.A. de C.V.* del estado de *Zacatecas*:

```

1  falta_producto('hotel del fresno de fresnillo s.a. de c.v.', 'zacatecas')

```

Resultado:

```

[{'p': {'estado': 'zacatecas',
        'municipio': 'fresnillo',
        'nombre_producto': 'botella de mezcal 1 lt.'}},
 {'p': {'estado': 'zacatecas',
        'municipio': 'fresnillo',
        'nombre_producto': 'botellas con mezcal 1lt.'}},
 {'p': {'estado': 'zacatecas',
        'municipio': 'fresnillo',
        'nombre_producto': 'botella de tequila reposado 100% de agave 950ml.'}},
 {'p': {'estado': 'zacatecas',
        'municipio': 'fresnillo',
        'nombre_producto': 'botella de tequila 950ml.'}}]

```

Figure 4: Resultado de la consulta 2.

Consulta 3: Dado un estado y un producto, buscar alternativas sin incumplimiento de ese producto o categoría.

Para la realización de esta consulta filtraremos a través de la sentencia WHERE para obtener los productos alternativos sin incumplimiento.

```

1  def buscar_alternativas(estado, producto):
2      query = f'''
3      MATCH (p:Producto)-[et:en_tienda]->(t:Tienda)
4      MATCH (e:Estado)

```

```

5      MATCH (m:Municipio)
6      MATCH (t)-->(m)-->(e)
7      WHERE et.incumplimiento <> "no se detecto incumplimiento"
8            AND e.nombre_estado = "{estado}" AND p.nombre_producto = "{producto}"
9      RETURN p,t,m,e
10     '''
11     alternativa = session.run(query)
12     return alternativa.data()

```

```

1 nombre_producto = 'cajacon reproductor de disco...'
2 alternativas = buscar_alternativas('mexico', nombre_producto)
3 print('Alternativas de lugares...')
4 for i, alternativa in enumerate(alternativas):
5     print(f"Alternativa {i+1}: {alternativa['m']['nombre_municipio']}, \
6           {alternativa['t']['razon_social']}")

```

Salida:

```

--Alternativas de lugares donde se puede encontrar sin incumplimiento el producto--
Alternativa 1: valle de chalco solidaridad, tiendas soriana s.a. de c.v.
Alternativa 2: cuautitlan izcalli, tiendas soriana s.a. de c.v.
Alternativa 3: naucalpan de juarez, tiendas soriana s.a. de c.v.

```

Consulta 4: Llevar un registro de las compras que ha hecho cada usuario y el lugar donde las ha hecho.

Para esta consulta se crea una función que permite agregar usuarios que compraron un producto. Se almacena la hora de la compra para distinguirla de haber una segunda compra del producto.

```

1 def registrar_compra(user,curp,tienda,producto,municipio,año,mes,dia):
2     q_compra_usuario=f'''
3     MERGE (p: Producto {{nombre_producto:"{producto}",municipio:"{municipio}"
4     ,estado:"{estado}"}})
5     MERGE (t: Tienda {{razon_social:"{tienda}"}})
6     MERGE (m: Municipio {{nombre_municipio:"{municipio}"}})
7     MERGE (p)-[:en_tienda]->(t)
8     MERGE (t)-[:en_municipio]->(m)
9     MERGE (u: Usuario {{CURP:"{curp}",nombre:"{user}"}})
10    MERGE (u)-[cp:compro_producto
11    {{hora:date({{year: {año}, month: {mes}, day: {dia}})}}]->(p)
12
13    RETURN u,p,t,m,cp.hora
14    '''
15
16    return session.run(q_compra_usuario).data()

```

Después usamos la siguiente función para regresar el historial de compra de un usuario

```

1 def historial_usuario(user,curp):
2
3     q_historial=f'''
4     MATCH (u: Usuario {{CURP:"{curp}",nombre:"{user}"}})
5     MATCH (u)-[cp:compro_producto]->(p:Producto)
6     RETURN u,p,cp.hora
7     '''
8
9     l=session.run(q_historial).data()
10
11     municipio=[]

```

```

11 producto=[]
12 estado=[]
13 hora=[]
14 for li in l:
15     municipio.append(li['p']['municipio'])
16     estado.append(li['p']['estado'])
17     producto.append(li['p']['nombre_producto'])
18     hora.append(li['cp.hora'])
19
20 return pd.DataFrame({'user':user,'estado':estado,'municipio':municipio,
21                      'hora':hora,'producto': producto})

```

Después de hacer algunas inserciones usando la función *registrar_compra* obtenemos la siguiente salida para el usuario *K. Reeves*

	user	estado	municipio	hora	producto
0	K. Reeves	tamaulipas	reynosa	2021-11-21	atun aleta amarilla en aceite marca tuny cont....
1	K. Reeves	chihuahua	chihuahua	2019-07-23	25 lapices delineadores
2	K. Reeves	hidalgo	pachuca de soto	2021-03-25	palomitas para microondas bolsa con 85 gramos ...
3	K. Reeves	aguascalientes	aguascalientes	2019-06-20	manos libres 1 pieza color blanco marca flet m...
4	K. Reeves	mexico	tecamac	2020-02-28	tarro grande comercializadora bofian

Figure 5: Historial de compras de K.Reeves

de esta manera podemos hacer la consulta, obtener el historial de compras y después buscar información como el patrón de compra de la persona.

Podemos ver el efecto que esto tuvo en la base de datos en la siguiente imagen:

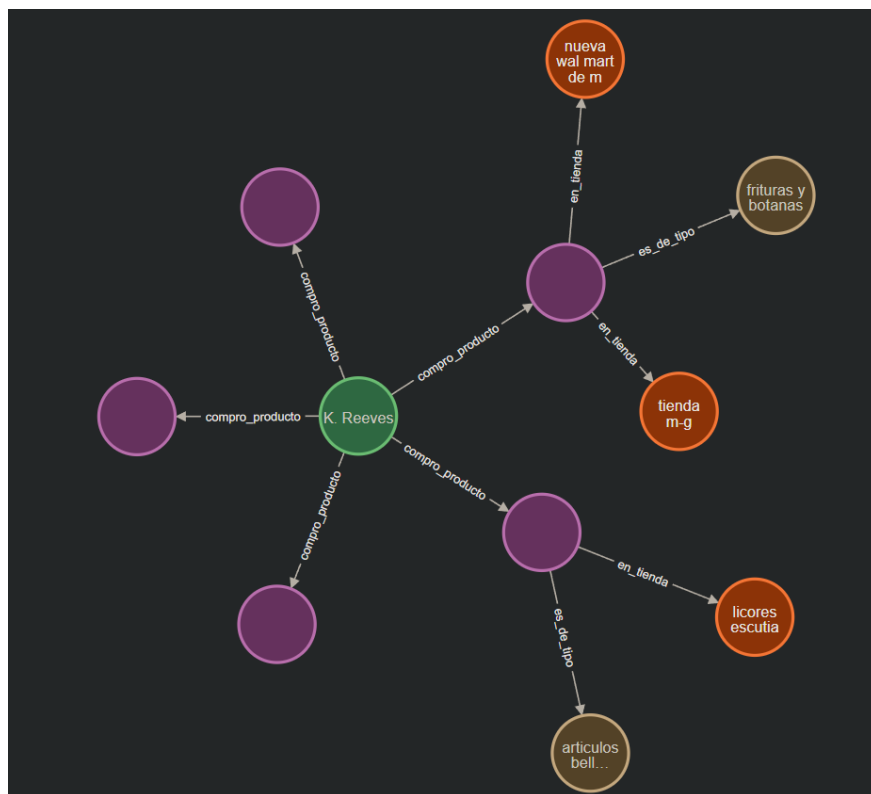


Figure 6: Muestra de la inserción de artículos comprados por *K. Reeves*. En la interfaz gráfica se desplegaron la tienda y el tipo de artículo de los dos artículos de la derecha para mostrar que las compras se unen al resto de la base.

Consulta 5: Encontrar los estados con mayor y menor incumplimiento relativo al número de tiendas que tiene.

La estrategia para resolver este problema será mediante la creación de dos consultas a la base de datos de Neo4j para obtener:

- Los productos y tiendas que presentan incumplimiento dado un estado.
- La cuenta del número total de tiendas en un estado.

El resultado de estas consultas es procesado para obtener la razón entre el número de tiendas que presentan incumplimiento y el número total de tiendas.

```

1  def incumplimiento_por_tienda(estado):
2
3      query = f"""
4      MATCH (e:Estado {{nombre_estado:"{estado}"}})
5      MATCH (m:Municipio)
6      MATCH (p:Producto {{estado:"{estado}", municipio:m.nombre_municipio}})
7      MATCH (t:Tienda)
8      MATCH (t) -[:en_municipio]-> (m)
9      MATCH (m) -[:en_estado]-> (e)
10     MATCH (p) -[:et:en_tienda]-> (t)
11     WHERE NOT et.incumplimiento = "no se detecto incumplimiento"
12     RETURN p, t
13     """
14
15     return session.run(query).data()
16
17 def tiendas_por_estado(estado):
18
19     query = f"""
20     MATCH (e:Estado {{nombre_estado:"{estado}"}})
21     MATCH (m:Municipio)
22     MATCH (p:Producto {{estado:"{estado}", municipio:m.nombre_municipio}})
23     MATCH (t:Tienda)
24     MATCH (t) -[:en_municipio]-> (m)
25     MATCH (m) -[:en_estado]-> (e)
26     MATCH (p) -[:et:en_tienda]-> (t)
27     RETURN t
28     """
29
30     return session.run(query).data()
31
32 faltantes = {}
33 for estado in profeco_df.ENTIDAD.unique():
34
35     tiendas = set()
36     for dict in incumplimiento_por_tienda(f'{estado}'):
37         tiendas.add(dict['t']['razon_social'])
38
39     tiendas_totales = set()
40     for dict in tiendas_por_estado(f'{estado}'):
41         tiendas_totales.add(dict['t']['razon_social'])
42
43     faltantes[f'{estado}'] = len(tiendas) / len(tiendas_totales)
44
45 # creamos un DataFrame con los resultados obtenidos
46 faltas_por_estado = pd.Series(faltantes).sort_values(ascending=False)
47

```

Graficando los resultados guardados en el DataFrame, se tiene la siguiente gráfica:

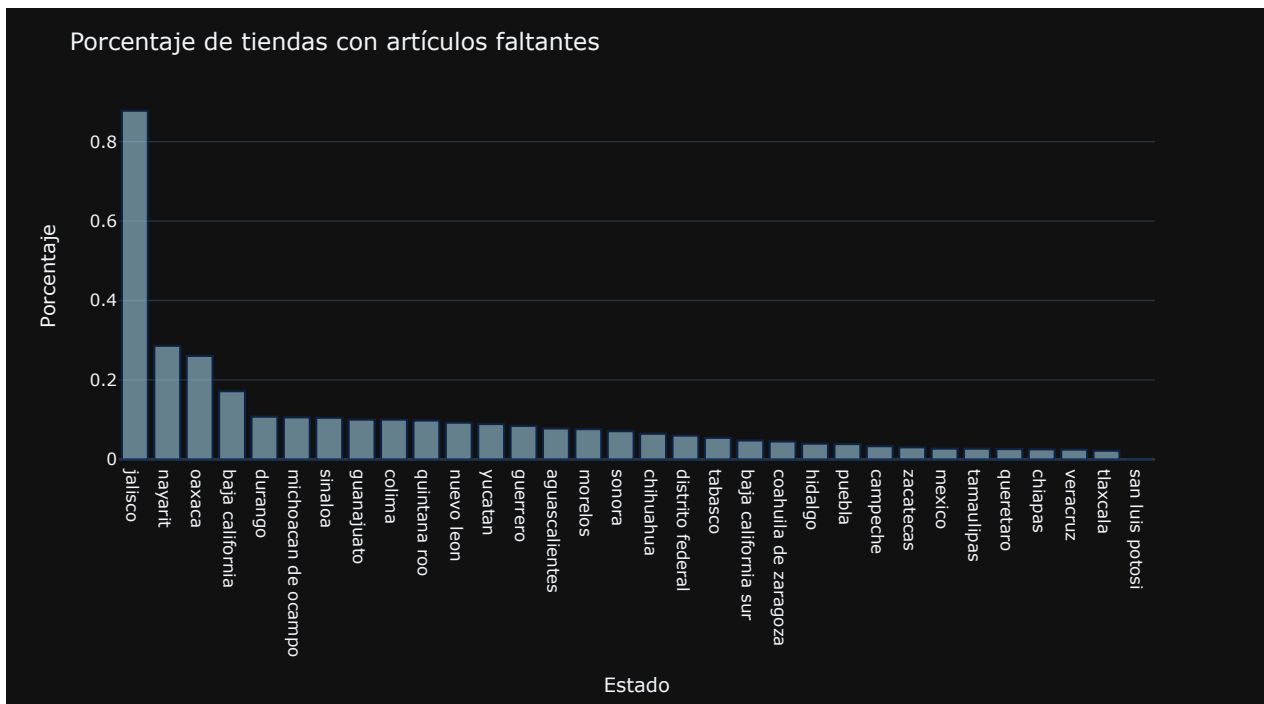


Figure 7: Resultado de la consulta 5.

Jalisco es el estado que la mayor de las irregularidades.

Conclusiones

Las bases de datos orientadas a grafos ofrecen una forma visual y sencilla de visualización de las relaciones de diferentes entidades, con lo cual es más sencillo el diseñar e implementar las bases de datos.

Estas bases de datos son particularmente buenas con las relaciones que existen entre diferentes objetos, lo cual la convierte en una base de datos ideal para el manejo de bases de datos con redes colaborativas. La sintaxis y acciones que se pueden ejecutar son bastante similares a las ya vistas en otras bases de datos, por lo que migrar a esta es un paso natural para el entendimiento de estructuras de datos más versátiles.