

Universidad Complutense de Madrid
Máster Big Data & Data Engineering

Bases de Datos NoSQL

Actividad Práctica

Profesor:
Marlon Cárdenas

Andrés Vega Zamora

13 de febrero de 2025

Índice

| | |
|---|-----------|
| 1. Exploración inicial de los datos | 3 |
| 1.1. Revisión inicial | 3 |
| 1.2. Revisión de campos | 3 |
| 1.3. Relación entre ficheros | 4 |
| 2. Modelado de datos | 5 |
| 2.1. Diseño del modelo de datos | 5 |
| 2.2. Validación del modelo mediante consultas | 7 |
| 2.3. Creación y uso de índices en MongoDB | 14 |
| 2.4. Modelo de datos: Extensión con alojamientos turísticos | 17 |
| 2.4.1. Consultas para evaluar el modelo | 20 |
| 3. Visualización de datos: Modelo de grafo | 29 |
| 3.1. Creación del grafo | 29 |
| 3.2. Consultas de validación | 31 |

1. Exploración inicial de los datos

1.1. Revisión inicial

Fichero Actividad económica y Locales:

Provenientes del Censo de Actividades y Locales de la Comunidad de Madrid. Se les asocian características físicas de los locales, actividades realizadas en el mismo y sus titulares.

Se cuenta con dos ficheros:

- Fichero de locales: Cada registro corresponde a un local.
- Fichero de actividades: Cada registro corresponde a un local más un código de actividad. Un local se repite tantas veces como actividades asociadas tenga.

Fichero Licencias:

Provee información física de los locales además de información sobre las licencias a concedidas a cada local.

Fichero Terrazas:

Provee información física de las terrazas y los locales a los que pertenecen.

1.2. Revisión de campos

Se realiza una limpieza de datos en todos los ficheros para evitar datos que no brindan información útil y datos repetidos.

Fichero locales

Se eliminan columnas que tienen muy pocos datos no nulos ("hora_apertura2", "hora_cierre2"):

```
1 df_locales = df_locales.drop(['hora_apertura2', 'hora_cierre2'], axis = 1, )
```

Se eliminan columnas que tienen un único valor en toda la columna:

```
1 df_locales= df_locales.drop(['fx_datos_ini', 'fx_datos_fin', 'id_clase_ndp_edificio', 'id_clase_ndp_acceso'], axis = 1);
```

Se eliminan los leading y trailing whitespaces de las columnas de texto:

```
1 for column in df_locales.columns:
2     if df_locales[column].dtype == "object":
3         df_locales[column] = df_locales[column].str.strip()
```

Se modifica el formato de ubicación geográfica a latitud y longitud:

```
1 from pyproj import CRS, Transformer
2
3 # 1. Define la proyeccion UTM (zona 30N para Madrid) y la proyeccion geografica (
4   WGS84)
5 utm_crs = CRS("EPSG:25830") # UTM Zone 30N
6 wgs84_crs = CRS("EPSG:4326") # WGS84 (latitud y longitud)
```

```

7 # 2. Crea el transformador de coordenadas
8 transformer = Transformer.from_crs(utm_crs, wgs84_crs)
9
10 # 3. Funcion para convertir coordenadas UTM a latitud y longitud
11 def utm_to_latlon(x, y):
12     lat, lon = transformer.transform(x, y)
13     return lat, lon
14
15 df_locales[['latitud', 'longitud']] = df_locales.apply(lambda row: utm_to_latlon(
    row['coordenada_x_local'], row['coordenada_y_local']), axis=1, result_type='
    expand')

```

Fichero actividades

Se convierte el tipo de columnas 42 y 44 a int64, ya que todos menos uno de sus valores son numéricos:

```

1 df_actividades['id_division'] = pd.to_numeric(df_actividades['id_division'],
    errors = 'coerce')
2 df_actividades['id_epigrafe'] = pd.to_numeric(df_actividades['id_epigrafe'],
    errors = 'coerce')
3 df_actividades['id_division'] = df_actividades['id_division'].replace(np.nan,
    999).infer_objects(copy=False)
4 df_actividades['id_epigrafe'] = df_actividades['id_epigrafe'].replace(np.nan,
    999).infer_objects(copy=False)
5 df_actividades['id_division'] = df_actividades['id_division'].astype('Int64')
6 df_actividades['id_epigrafe'] = df_actividades['id_epigrafe'].astype('Int64')

```

Se nota que cuando el id_epigrafe es -1, no hay informacion del local. Se eliminan esas filas:

```

1 df_actividades = df_actividades[df_actividades["id_epigrafe"] != -1]

```

Fichero terrazas

Se eliminan los leading y trailing whitespaces de las columnas de texto:

```

1 for column in df_locales.columns:
2     if df_locales[column].dtype == "object":
3         df_locales[column] = df_locales[column].str.strip()

```

Fichero licencias

Se cambia el formato de la fecha a ISO Date YYYY-MM-DD:

```

1 df_licencias["Fecha_Dec_Lic"] = pd.to_datetime(df_licencias["Fecha_Dec_Lic"],
    format = '%d/%m/%Y')
2 df_licencias["Fecha_Dec_Lic"] = df_licencias["Fecha_Dec_Lic"].dt.strftime('%Y-%m
    -%d')

```

1.3. Relación entre ficheros

Luego de observar los datos pertenecientes a cada fichero se nota que todos los ficheros tienen muchas columnas que se repiten con respecto al fichero de locales, y estas columnas dan información física

sobre los locales en sí, no de las terrazas, licencias o actividades. De manera que se procede a eliminar las columnas repetidas.

```
1 #Se eliminan columnas repetidas entre fichero locales y actividades, del fichero
  de actividades
2
3 columns_locales_actividades = set(df_locales.columns) & set(df_actividades.
  columns)
4 columns_locales_actividades.remove('id_local')
5 df_actividades = df_actividades.drop(columns_locales_actividades, axis = 1)
```

Se sigue el mismo proceso para las columnas entre locales-terrazas y locales-licencias. De manera cada uno de estos ficheros ahora solo brinda información sobre su propio ámbito, y la relación con el fichero de locales se establece por medio de la columna “id_local” que mantienen todos los ficheros.

2. Modelado de datos

2.1. Diseño del modelo de datos

Se decide utilizar un modelo desnormalizado/embebido en los locales para ingresar los datos en MongoDB. Se prioriza la eficiencia en la consulta de datos sobre la eficiencia en la actualización de datos.

Se establece la conexión con Mongo:

```
1 client = MongoClient("mongodb://localhost:27017/")
2 db = client["locales_madrid_mongo"]
3 locales_collection = db["locales_madrid"]
```

Se convierten los dataframes a listas de diccionarios:

```
1 datos_locales = df_locales.to_dict(orient = "records")
2 datos_actividades = df_actividades.to_dict(orient = "records")
3 datos_licencias = df_licencias.to_dict(orient = "records")
4 datos_terrazas = df_terrazas.to_dict(orient = "records")
```

Construcción del modelo embebido:

```
1 datos_embebidos = {}
2
3 for actividad in datos_actividades:
4     id_local = actividad["id_local"]
5     del actividad["id_local"] # Se elimina el id_local del diccionario de
  actividad porque ya la actividad va a estar embebida dentro del local
6     datos_embebidos.setdefault(id_local, {"actividades" : [], "licencias" : [], "
  terraza": []}) #Vacio si no existen los datos
7     datos_embebidos[id_local]["actividades"].append(actividad)
8
```

```

9 for licencia in datos_licencias:
10     id_local = licencia["id_local"]
11     del licencia["id_local"]
12     datos_embebidos.setdefault(id_local, {"actividades" : [], "licencias" : [], "
terrazza": []}) #Vacio si no existen los datos
13     datos_embebidos[id_local]["licencias"].append(licencia)
14
15 for terraza in datos_terrazas:
16     id_local = terraza["id_local"]
17     del terraza["id_local"]
18     datos_embebidos.setdefault(id_local, {"actividades" : [], "licencias" : [], "
terrazza": []}) #Vacio si no existen los datos
19     datos_embebidos[id_local]["terrazza"].append(terrazza)

```

Inserción de datos a Mongo:

```

1 for local in datos_locales:
2     id_local = local["id_local"]
3     info_embebida = datos_embebidos.get(id_local)
4     if info_embebida:
5         local.update(info_embebida)
6     locales_collection.insert_one(local)
7 client.close()

```

La colección resultante tiene la siguiente estructura:

```

{
  "_id": {"$oid": "67a7ad7ab7911aea08fe75b7"},
  "id_local": 40002093,
  "id_distrito_local": 4,
  "desc_distrito_local": "SALAMANCA",
  "id_barrio_local": 406,
  "desc_barrio_local": "CASTELLANA",
  ...
  "actividades": [
    {
      "id_seccion": "I",
      "desc_seccion": "HOSTELERÍA",
      "id_division": 56,
      "desc_division": "SERVICIOS DE COMIDAS Y BEBIDAS",
      "id_epigrafe": 561005,
      "desc_epigrafe": "BAR CON COCINA"
    }
  ],
  "licencias": [

```

```

{
  "ref_licencia": "350/2023/36134",
  "id_tipo_licencia": 4,
  "desc_tipo_licencia": "Transmisión de licencia Urbanística",
  "id_tipo_situacion_licencia": 5,
  "desc_tipo_situacion_licencia": "Transmisión de Licencia Concedida",
  "Fecha_Dec_Lic": "1900-01-01"
}
],
"terraza": [{...}]
}

```

2.2. Validación del modelo mediante consultas

Se ejecutan las siguientes consultas para evaluar la eficiencia y funcionalidad del modelo. Se adjuntan **recortes** de los resultados de las consultas.

a - Total de locales y terrazas por distrito y barrio

```

1  [
2    {
3      $group: { //Agrupar por barrio y distrito
4        _id: {
5          barrio: "$desc_barrio_local",
6          distrito: "$desc_distrito_local"
7        },
8        totalLocales: { $sum: 1 },
9        totalTerrazas: { //Verifica si el local tiene terraza
10         $sum: {
11           $cond: {
12             if: {
13               $and: [
14                 { $not: { $eq: ["$terraza", []] } },
15                 { $gt: [{ $size: { $ifNull: ["$terraza", []] } }, 0] }
16               ]
17             },
18             then: 1,
19             else: 0
20           }
21         }
22       }
23     },
24   ],
25   {
26     $project: { //Selecciona campos a mostrar

```

```

27     _id: 0, // Excluye el campo _id
28     distrito: "$_id.distrito",
29     barrio: "$_id.barrio",
30     totalLocales: 1,
31     totalTerrazas: 1
32   }
33 },
34 {
35   $sort: { // Ordena por distrito y barrio
36     distrito: 1,
37     barrio: 1
38   }
39 }
40 ]

```

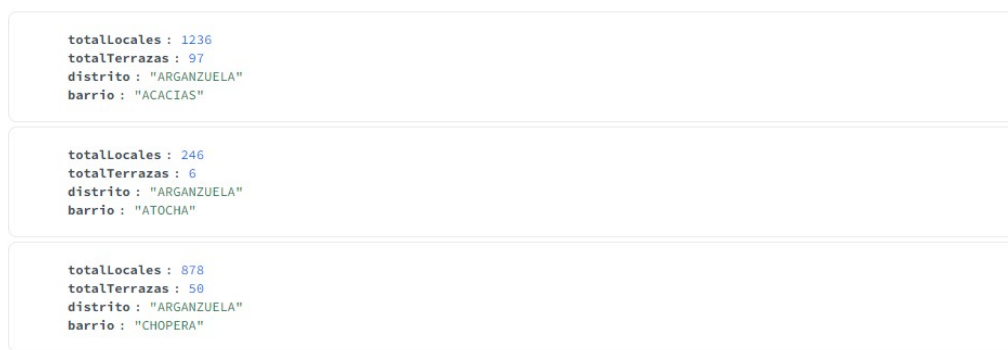


Figura 1: Locales y terrazas por distrito y barrio

b - Tipos de licencias y cantidad de licencias por tipo

```

1 [
2   {
3     $unwind: "$licencias" // Desglosa el array de licencias
4   },
5   {
6     $group: {
7       _id: "$licencias.desc_tipo_licencia", // Agrupa por tipo de licencia
8       cantidad: { $sum: 1 } // Cuenta las licencias de cada tipo
9     }
10  },
11  {
12    $project: { //Selecciona los campos a mostrar
13      _id: 0, // Excluye el campo _id
14      tipo_licencia: "$_id",
15      cantidad: 1
16    }
17  },

```



```

18 {
19   $sort: { // Ordena por cantidad de licencias de mayor a menor
20     cantidad: -1
21   }
22 }
23 ]

```

| | |
|------------------|--|
| cantidad : 57028 | tipo_licencia : "Transmisión de licencia Urbanística" |
| cantidad : 51249 | tipo_licencia : "Declaración Responsable" |
| cantidad : 30783 | tipo_licencia : "Licencia Urbanística" |
| cantidad : 5956 | tipo_licencia : "Licencia recogida en el trabajo de campo" |

Figura 2: Tipos de licencias y cantidad

c - Listado de locales y terrazas con licencias “En tramitación”

```

1 [
2   {
3     $unwind: "$licencias" // Desglosa el array de licencias
4   },
5   {
6     $match: { // Verifica que la licencia este "En tramitacion"
7       $expr: {
8         $regexMatch: {
9           input:
10             "$licencias.desc_tipo_situacion_licencia",
11           regex: "^En tramitacion$",
12           options: "i"
13         }
14       }
15     }
16   },
17   {
18     $group: { //Agrupar los resultados
19       _id: "$id_local",
20       rotulo: { $first: "$rotulo" },
21       licencias_en_tramitacion: { $sum: 1 },
22       terraza: { $first: "$terrazza" }
23     }
24   },
25   {
26     $project: { //Selecciona los campos a mostrar

```

```

27     _id: 0,
28     id_local: "$_id",
29     rotulo: 1,
30     terraza: 1,
31     licencias_en_tramitacion: 1
32   }
33 }
34 ]

```

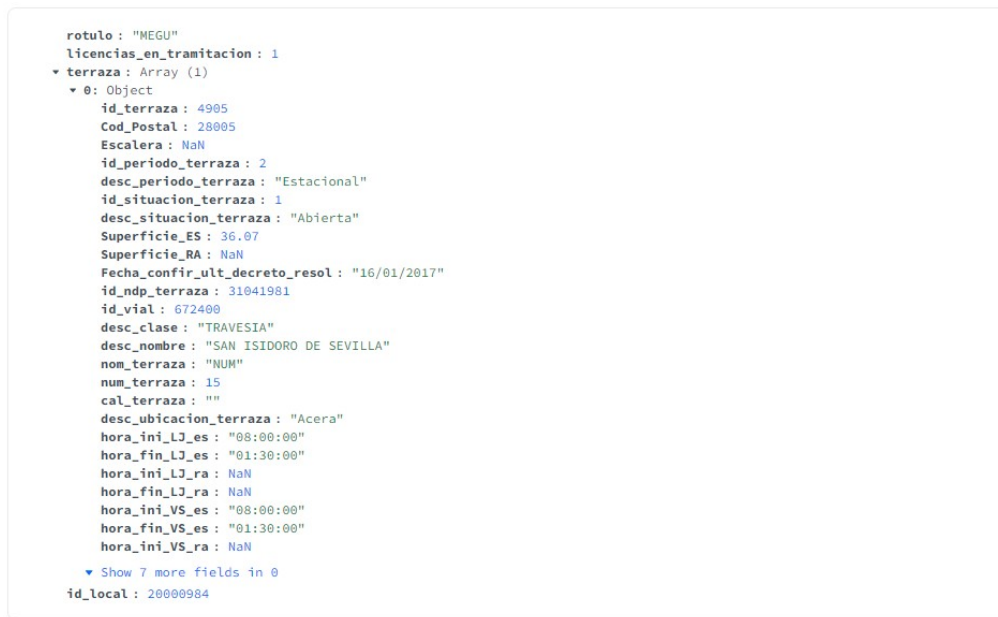


Figura 3: Locales y terrazas “En tramitación”

d - Clasificación por sección, división y epígrafe

```

1  [
2    {
3      $unwind: "$actividades" // Se desglosa el array de actividades
4    },
5    {
6      $group: { //Se agrupa por seccion, division y epigrafe, y se maneja el
7        caso vacio
8        _id: {
9          seccion: { $ifNull: ["$actividades.id_seccion", "Sin seccion"] },
10         division: { $ifNull: ["$actividades.id_division", "Sin division"] },
11         epigrafe: { $ifNull: ["$actividades.id_epigrafe", "Sin epigrafe"] },
12         //
13         desc_epigrafe: { $ifNull: ["$actividades.desc_epigrafe", "Sin
14           describir epigrafe"]}
15       },
16       totalLocales: { $sum: 1 }, // Se cuentan los locales en cada grupo
17     }
18   ]

```

```

14     locales: {
15         $push: { // Se agregan los locales a un array dentro de cada grupo
16             id_local: "$id_local",
17             rotulo: "$rotulo",
18             terraza: "$terrazza"
19         }
20     }
21 }
22 },
23 {
24     $project: { // Se seleccionan los campos a mostrar
25         _id: 0,
26         seccion: "$_id.seccion",
27         division: "$_id.division",
28         epigrafe: "$_id.epigrafe",
29         totalLocales: 1,
30         locales: 1,
31         desc_epigrafe: "$_id.desc_epigrafe"
32     }
33 },
34 {
35     $sort: { // Se ordena por total de locales
36         totalLocales:-1
37     }
38 }
39 ]

```

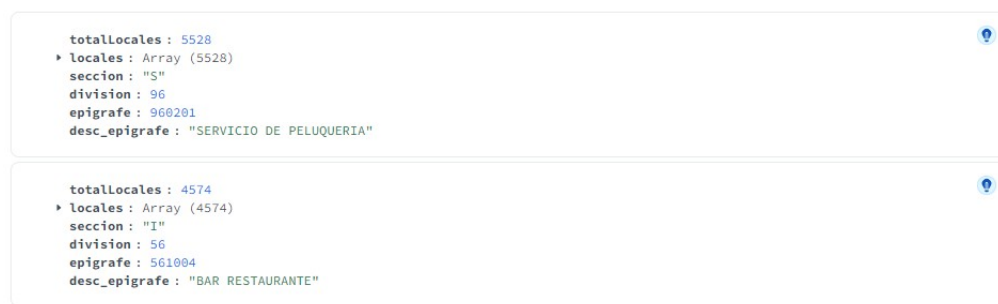


Figura 4: Clasificación de actividades

e - Actividad más frecuente por distrito y barrio

```

1 [
2   {
3     $unwind: "$actividades" // Se desglosa el array de actividades
4   },
5   {
6     $group: { //Se agrupa por distrito, barrio y actividad

```

```
7      _id: {
8        desc_distrito_local: "$desc_distrito_local",
9        desc_barrio_local: "$desc_barrio_local",
10       desc_epigrafe: "$actividades.desc_epigrafe"
11     },
12     count: { $sum: 1 } // Cuenta la frecuencia de cada combinacion de
13     distrito, barrio y actividad
14   },
15   {
16     $sort: { count: -1 } // Ordena los resultados por frecuencia de forma
17     descendente
18   },
19   {
20     $group: { // Se agrupa solo por distrito y barrio
21       _id: {
22         desc_distrito_local: "$_id.desc_distrito_local",
23         desc_barrio_local: "$_id.desc_barrio_local"
24       },
25       actMasFrecuente: { $first: "$_id.desc_epigrafe" }, // Obtiene la
26       actividad mas frecuente
27       cuentaMasFrecuente: { $first: "$count" } // Obtiene la frecuencia de
28       la actividad mas frecuente
29     }
30   },
31   {
32     $project: { //Se seleccionan los campos a mostrar
33       _id: 0,
34       distrito: "$_id.desc_distrito_local",
35       barrio: "$_id.desc_barrio_local",
36       actividad_mas_frecuente: "$actMasFrecuente",
37       frecuencia: "$cuentaMasFrecuente"
38     }
39   },
40   {
41     $sort: {
42       frecuencia: -1
43     }
44   }
45 ]
```

| |
|---|
| distrito : "SALAMANCA" barrio : "RECOLETOS" actividad_mas_frecuente : "COMERCIO AL POR MENOR DE PRENDAS DE VESTIR EN ESTABLECIMIENTOS ESPECIALIZADOS" frecuencia : 329 |
| distrito : "VILLA DE VALLECAS" barrio : "CASCO H.VALLECAS" actividad_mas_frecuente : "COMERCIO AL POR MAYOR FRUTAS, VERDURAS Y DERIVADOS" frecuencia : 216 |
| distrito : "CENTRO" barrio : "UNIVERSIDAD" actividad_mas_frecuente : "COMERCIO AL POR MENOR DE PRENDAS DE VESTIR EN ESTABLECIMIENTOS ESPECIALIZADOS" frecuencia : 291 |

Figura 5: Actividad más frecuente por barrio

f - Cambio de horarios de Bares con Cocina en SALAMANCA

```

1 //filtrado
2 {
3   "actividades.desc_epigrafe": "BAR CON COCINA",
4   "desc_distrito_local": "SALAMANCA"
5 }

```

```

1 //actualizacion
2 {
3   $set: {
4     hora_apertura1: "11:00",
5     hora_cierre1: "01:00"
6   }
7 }

```

```

1 //Consulta para verificar cambio
2 [
3   {
4     $match: {
5       "actividades.desc_epigrafe": "BAR CON COCINA",
6       "desc_distrito_local": "SALAMANCA" // Filtrar por el distrito de
7       Salamanca
8     }
9   },
10  {
11    $unwind: "$terrazas" // Desglosa el array de terrazas para acceder a los
12    horarios
13  },
14  {
15    $project: { // Selecciona los campos a mostrar
16      _id: 0,
17      nombre_local: "$rotulo",
18      distrito: "$desc_distrito_local",
19      hora_apertura: "$hora_apertura1",

```

```
18     hora_cierre: "$hora_cierre1"
19   }
20 }
21 ]
```

2.3. Creación y uso de índices en MongoDB

Se crean tres tipos de índices distintos y luego se ejecutan consultas para evaluar su uso y desempeño:

a - Índice simple: Nombre del barrio

Creación del índice:

```
1 db.locales_madrid.createIndex({desc_barrio_local: 1})
```

Consulta para evaluar resultados:

```
1 //Consulta de locales con terraza en barrio que empiece con "A"
2 [
3   {
4     $match: {
5       "terrazza": { $exists: true, $not: { $size: 0 } }, // Filtra locales
6       "desc_barrio_local": { $regex: "^A" } // Filtra barrios que empiezan
7       con "A"
8     },
9     {
10      $group: {
11        _id: "$desc_barrio_local", // Agrupa por barrio
12        totalLocalesConTerraza: { $sum: 1 } // Cuenta locales con terraza en
13        cada barrio
14      },
15      {
16        $sort: { _id: 1 } // Ordena alfabeticamente por nombre de barrio
17      },
18      {
19        $project: { // Selecciona los campos a mostrar
20          _id: 0,
21          nombre_barrio: "$_id",
22          totalTerrazas: "$totalLocalesConTerraza"
23        }
24      }
25    ]
```

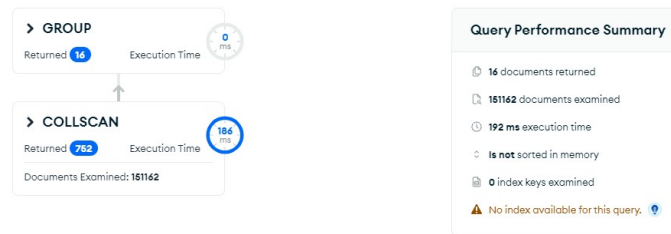


Figura 6: Rendimiento de la consulta sin el index simple

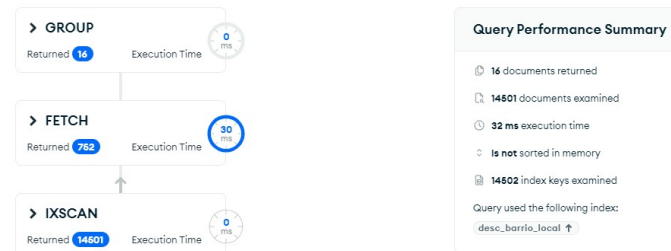


Figura 7: Rendimiento de la consulta con el index simple

Al comparar las figuras 6 y 7 se puede apreciar que el tiempo de la consulta baja de 192ms a 32 ms al usar el índice simple para desc_barrio_local.

b - Índice compuesto: Distrito y barrio Creación del índice:

```
1 db.locales_madrid.createIndex({desc_distrito_local: 1, desc_barrio_local:
  1})
```

Consulta para evaluar resultados:

```
1 //Consulta de locales con terraza en distrito que empiece con C y Barrio A
2 [
3   {
4     $match: {
5       terraza: {
6         $exists: true,
7         $not: { $size: 0 }
8       }, // Filtra locales con terraza
9       desc_distrito_local: { $regex: "^C" }, // Filtra distritos que
        empiezan con "C"
10      desc_barrio_local: { $regex: "^A" } // Filtra barrios que empiezan con
        "A"
11    }
12  },
13  {
14    $group: { //Agrupar por distrito y barrio
15      _id: {
```

```

16     distrito: "$desc_distrito_local",
17     barrio: "$desc_barrio_local"
18   },
19   totalLocalesConTerraza: { $sum: 1 }
20 }
21 },
22 {
23   $sort: { "_id.distrito": 1, "_id.barrio": 1 } // Ordena por distrito y
24   luego por barrio
25 },
26 {
27   $project: { // Selecciona los campos a mostrar
28     _id: 0,
29     distrito: "$_id.distrito",
30     barrio: "$_id.barrio",
31     totalTerrazas: "$totalLocalesConTerraza"
32   }
33 ]

```

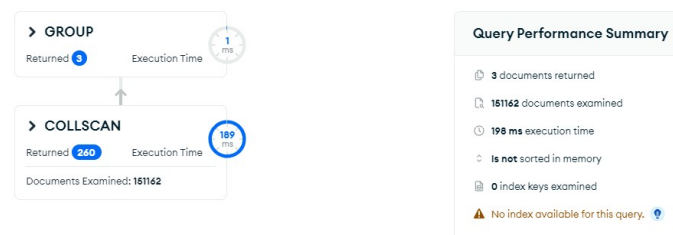


Figura 8: Rendimiento de la consulta sin el index compuesto

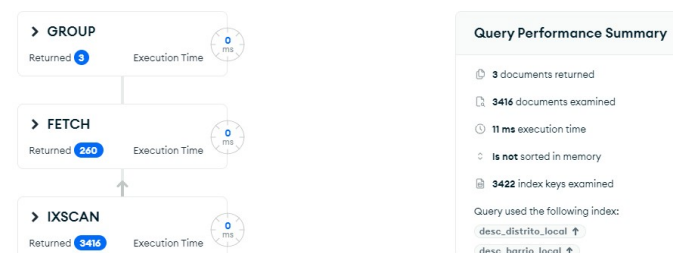


Figura 9: Rendimiento de la consulta con el index compuesto

Al comparar las figuras 8 y 9 se puede apreciar que el tiempo de la consulta baja de 198ms a 11 ms al usar el índice compuesto para desc_distrito_local y desc_barrio_local.

c - Índice de array: Descripción de actividad Creación del índice:

```

1 db.locales_madrid.createIndex({"actividades.desc_epigrafe": 1})

```


Consulta para evaluar resultados:

```

1 //Consulta de locales donde al menos una de sus actividades sea BAR CON
  COCINA
2 [
3   {
4     $match: {
5       "actividades.desc_epigrafe": "BAR CON COCINA"
6     }
7   },
8   {
9     $count: "totalLocalesConBarCocina"
10  }
11 ]

```

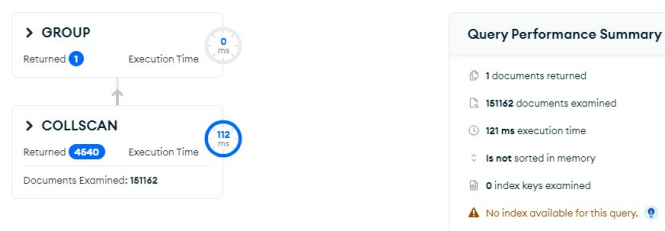


Figura 10: Rendimiento de la consulta sin el index de array

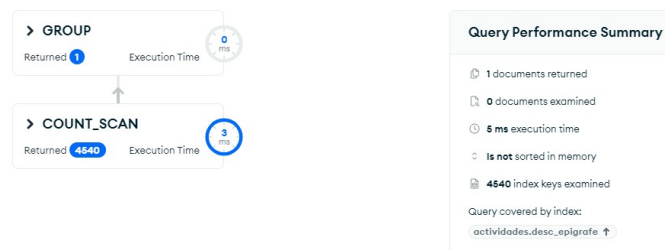


Figura 11: Rendimiento de la consulta con el index de array

Al comparar las figuras 10 y 11 se puede apreciar que el tiempo de la consulta baja de 121ms a 6 ms al usar el índice compuesto para desc_distrito_local y desc_barrio_local.

2.4. Modelo de datos: Extensión con alojamientos turísticos

Se incorporan datos de alojamientos turísticos de Airbnb en Madrid. Tomados de <https://insideairbnb.com/madrid/>. El fichero tiene 75 columnas, sin embargo para las consultas de la tarea muchas no son necesarias.

```

1 columns_to_keep = ["id", "name", "neighbourhood_cleansed", "
    neighbourhood_group_cleansed", "latitude", "longitude", "bathrooms", "bedrooms
    ", "amenities", "price", "number_of_reviews", "review_scores_rating"]
2 df_airbnb = df_airbnb[columns_to_keep]

```

Además nos aseguramos de limpiar los trailing y leading whitespaces de las columnas de texto:

```

1 for column in df_airbnb.columns:
2     if df_airbnb[column].dtype == "object":
3         df_airbnb[column] = df_airbnb[column].str.strip()

```

Limpiamos el formato de la columna “price”. Eliminamos el “\$”, eliminamos la “,” y convertimos la columna a float.

```

1 df_airbnb["price"] = df_airbnb["price"].str.replace("$", "")
2 df_airbnb["price"] = df_airbnb["price"].str.replace(",", "")
3 df_airbnb["price"] = df_airbnb["price"].astype(float)

```

Eliminamos las tildes y pasamos a mayúsculas los campos de neighbourhood_cleansed y neighbourhood_group_cleansed. Además reemplazamos los espacios antes y después de los guiones:

```

1 df_airbnb["neighbourhood_cleansed"] = df_airbnb["neighbourhood_cleansed"].str.
    normalize('NFKD').str.encode('ascii', 'ignore').str.decode('utf-8').str.upper
    ()
2 df_airbnb["neighbourhood_group_cleansed"] = df_airbnb["
    neighbourhood_group_cleansed"].str.normalize('NFKD').str.encode('ascii', '
    ignore').str.decode('utf-8').str.upper
3 df_airbnb.neighbourhood_group_cleansed = df_airbnb.neighbourhood_group_cleansed.
    str.replace(" - ", "-")

```

Cambiamos los nombres de las columnas para facilitar la redacción y lectura de las consultas:

```

1 nombres_columnas = ["id_airbnb", "nombre_airbnb", "barrio_airbnb", "
    distrito_airbnb", "latitud", "longitud", "numero_banos",
2     "numero_habitaciones", "servicios", "precio", "numero_resenas
    ", "puntuacion"]
3 df_airbnb.columns = nombres_columnas

```

Para insertar los datos a Mongo se van a manejar como una colección aparte en la misma base de datos. Para establecer las consultas que usan ambas colecciones, la relación se suele establecer por medio de los barrios o distritos de Madrid, ya que los valores se comparten entre ambas colecciones. Los servicios de cada alojamiento se anidan en un array.

Carga de datos a Mongo:

```

1 from pymongo import MongoClient
2 import ast
3
4 client = MongoClient("mongodb://localhost:27017/")
5 db = client["locales_madrid_mongo"]

```

```

6 alojamientos_collection = db["alojamientos"]
7
8 def preparar_datos(row):
9     """Prepara un diccionario para insertar en MongoDB."""
10    documento = row.to_dict()
11
12    # Aseguramos que "servicios" sea un array
13    if isinstance(documento["servicios"], str):
14        documento["servicios"] = [s.strip() for s in documento["servicios"].strip("[]")
15                                   .replace('\\"', '\"').split(",")]
16    elif not isinstance(documento["servicios"], list): # Si no es ni string ni
17        lista, lo dejamos vacio
18        documento["servicios"] = []
19
20    return documento
21
22 # Convertimos el DataFrame a una lista de diccionarios
23 datos_para_insertar = [preparar_datos(row) for _, row in df_airbnb.iterrows()]
24
25 # Insertamos los documentos
26 if datos_para_insertar: # Verificamos que haya datos para insertar
27     alojamientos_collection.insert_many(datos_para_insertar)
28     print(f"Se insertaron {len(datos_para_insertar)} documentos en la coleccion.")
29 else:
30     print("No hay datos para insertar.")
31
32 # Cerramos la conexion
33 client.close()

```

La estructura de los documentos en la colección “alojamientos” es la siguiente:

```

{
  "_id": {
    "$oid": "67a5d21accc82f73cdbb67d4"
  },
  "id_airbnb": 62423,
  "nombre_airbnb": "MAGIC ARTISTIC HOUSE IN THE CENTER OF MADRID",
  "barrio_airbnb": "JUSTICIA",
  "distrito_airbnb": "CENTRO",
  "latitud": 40.41884,
  "longitud": -3.69655,
  "numero_baños": 1.5,
  "numero_habitaciones": 1,
  "servicios": [
    "Books and reading material",

```

```

    "First aid kit",
    "Kitchen",
    "Essentials",
    ...
  ],
  "precio": 69,
  "numero_reseñas": 219,
  "puntuacion": 4.64,
  "id_barrio_alojamiento": 104,
  "id_distrito_alojamiento": 1
}

```

2.4.1. Consultas para evaluar el modelo

a - Total de alojamientos, locales y terrazas por distrito y barrio

```

1  [
2    {
3      $group: { //Agrupa los alojamientos por distrito y barrio
4        _id: {
5          distrito: "$distrito_airbnb",
6          barrio: "$barrio_airbnb"
7        },
8        total_alojamientos: { $sum: 1 }
9      }
10   },
11
12   {
13     $lookup: { // Trae datos de locales_madrid usando el pipeline
14       from: "locales_madrid",
15       let: { distrito: "$_id.distrito", barrio: "$_id.barrio" },
16       pipeline: [
17         {
18           $match: { //Busca que concidan el distrito y barrio con los
19             alojamiento
20             $expr: {
21               $and: [
22                 { $eq: ["$desc_distrito_local", "$$distrito"] },
23                 { $eq: ["$desc_barrio_local", "$$barrio"] }
24               ]
25             }
26           },
27           {
28             $group: { //Cuenta locales y terrazas

```

```
29         _id: null,
30         total_locales: { $sum: 1 },
31         total_terrazas: {
32             $sum: {
33                 $cond: [
34                     { $gt: [{ $size: { $ifNull: ["$terrazas", []] } }, 0] },
35                     1,
36                     0
37                 ]
38             }
39         }
40     }
41 },
42 ],
43 as: "locales"
44 }
45 },
46 {
47     $unwind: "$locales" // Desglosa el array "locales"
48 },
49 {
50     $project: { // Selecciona los campos mostrar
51         _id: 0,
52         distrito: "$_id.distrito",
53         barrio: "$_id.barrio",
54         total_alojamientos: 1,
55         total_locales: "$locales.total_locales",
56         total_terrazas: "$locales.total_terrazas"
57     }
58 },
59 {
60     $sort:{ // Ordena por distrito y barrio
61         distrito: 1,
62         barrio: 1
63     }
64 }
65 ]
```

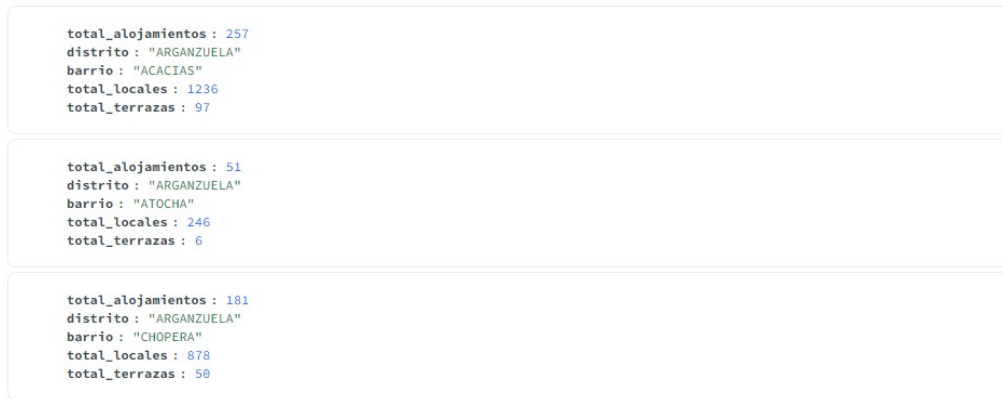


Figura 12: Alojamientos, locales y terrazas por barrio

b - 5 barrios con más alojamientos y terrazas con licencias concedidas hace 2 años o menos.

```

1 [
2   {
3     $group: {
4       _id: {
5         distrito: "$distrito_airbnb",
6         barrio: "$barrio_airbnb"
7       },
8       total_alojamientos: {
9         $sum: 1
10      }
11    }
12  },
13  {
14    $lookup: {
15      from: "locales_madrid",
16      let: {
17        distrito: "$_id.distrito",
18        barrio: "$_id.barrio"
19      },
20      pipeline: [
21        {
22          $match: {
23            $expr: {
24              $and: [
25                {
26                  $eq: [
27                    "$desc_distrito_local",
28                    "$$_id.distrito"
29                  ]
30                },
31                {
32                  $eq: [

```

```
33         "$desc_barrio_local",
34         "$$barrio"
35     ]
36 }
37 ]
38 }
39 }
40 },
41 {
42     $unwind: "$licencias"
43 },
44 {
45     $match: {
46         // Se agrega la condicion para desc_situacion_licencia
47         $expr: {
48             $and: [
49                 {
50                     $eq: [
51                         "$licencias.desc_tipo_situacion_licencia",
52                         "Concedida"
53                     ]
54                 },
55                 {
56                     $gte: [
57                         "$licencias.Fecha_Dec_Lic",
58                         "2023-02-08"
59                     ]
60                 }
61             ]
62         }
63     }
64 },
65 {
66     $group: {
67         _id: null,
68         total_terrazas: {
69             $sum: {
70                 $cond: [
71                     {
72                         $gt: [{ $size: { $ifNull: ["$terrazas", []] } }, 0] },
73                     1,
74                     0
75                 ]
76             }
77         }
78     }
79 }
```

```
80     ],
81     as: "locales"
82   }
83 },
84 {
85   $unwind: "$locales"
86 },
87 {
88   $project: {
89     _id: 0,
90     distrito: "$_id.distrito",
91     barrio: "$_id.barrio",
92     total_alojamientos: 1,
93     total_terrazas: "$locales.total_terrazas"
94   }
95 },
96 {
97   $addFields: {
98     suma_total: {
99       $add: [
100         "$total_alojamientos",
101         "$total_terrazas"
102       ]
103     }
104   }
105 },
106 {
107   $sort: {
108     suma_total: -1
109   }
110 },
111 {
112   $limit: 5
113 }
114 ]
115 ]
```

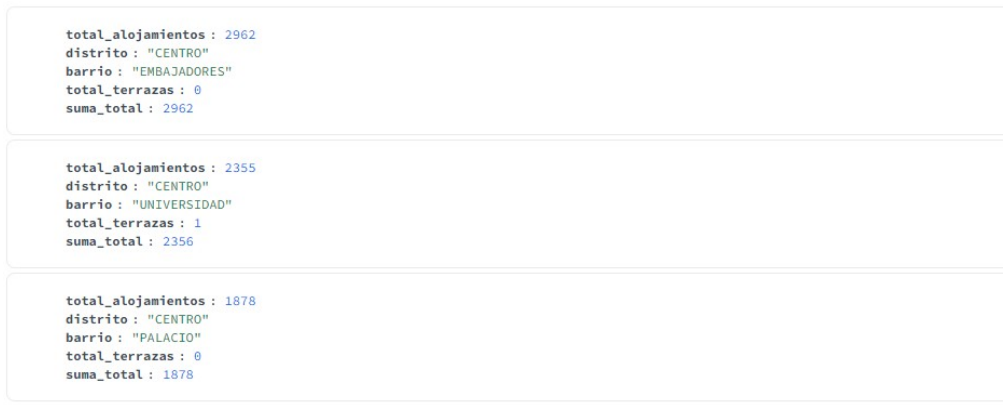



Figura 13: Cinco barrios con más alojamientos y terrazas concedidas hace menos de 2 años

c - Alojamientos con más reseñas por barrio y numero de locales con actividad comercial en dicho barrio

```

1 // Mas de 800 resenas
2 // Local que incluya COMERCIO AL POR MAYOR Y AL POR MENOR en desc_seccion
3
4 [
5   {
6     $match: // Se seleccionan alojamientos con mas de 800 resenas
7     {
8       numero_resenas: {
9         $gte: 800
10      }
11    }
12  },
13  {
14    $lookup: // Se traen los locales de los mismos barrios que los
15             alojamientos
16    {
17      from: "locales_madrid",
18      localField: "barrio_airbnb",
19      foreignField: "desc_barrio_local",
20      pipeline: [
21        {
22          $match: { // Se verifica que el local sea comercial
23            "actividades.desc_seccion": {
24              $regex:
25                /COMERCIO AL POR MAYOR Y AL POR MENOR/
26            }
27          }
28        },
29        as: "locales"

```

```

30     }
31 },
32 {
33     $project: // Se seleccionan los campos a mostrar
34     {
35         _id: 0,
36         nombre_airbnb: 1,
37         numero_resenas: 1,
38         locales_comerciales: {
39             $size: "$locales"
40         }
41     }
42 }
43 ]

```



Figura 14: Alojamientos con más reseñas por barrio, y locales comerciales en dicho barrio

d - Categorización de alojamientos por barrio, considerando precio, reseñas, numero de habitaciones y servicios

```

1 //precio
2 //ECONOMICO <= 70
3 //REGULAR >70 < 180
4 //CARO >= 180
5 //resenas (puntuacion)
6 //MALO <= 3
7 //ACEPTABLE >3 <4
8 //BUENO >= 4 <= 4,5
9 //EXCELENTE > 4,5
10 //numero de habitaciones
11 //PEQUENO 1
12 //MEDIANO 3 5
13 //GRANDE > 5
14 //servicios
15 //SIN SERVICIOS 0
16 //BASICO 10
17 //REGULAR 25
18 //MUCHOS SERVICIOS >25

```

```

19 [
20 {
21   $project: { // Se crean las categorias
22     precio_categoria: {
23       $switch: {
24         branches: [
25           { case: { $lte: ["$precio", 70] }, then: "ECONOMICO" },
26           { case: { $and: [{ $gt: ["$precio", 70] }, { $lte: ["$precio",
27 180] }] }, then: "REGULAR" }
28         ],
29         default: "CARO"
30       }
31     },
32     puntuacion_categoria: {
33       $switch: {
34         branches: [
35           { case: { $lt: ["$puntuacion", 3] }, then: "MALO" },
36           { case: { $and: [{ $gte: ["$puntuacion", 3] }, { $lte: ["
37 $puntuacion", 4] }] }, then: "ACEPTABLE" },
38           { case: { $and: [{ $gt: ["$puntuacion", 4] }, { $lte: ["
39 $puntuacion", 4.5] }] }, then: "BUENO" }
40         ],
41         default: "EXCELENTE"
42       }
43     },
44     habitaciones_categoria: {
45       $switch: {
46         branches: [
47           { case: { $eq: ["$numero_habitaciones", 1] }, then: "PEQUENO" },
48           { case: { $and: [{ $gte: ["$numero_habitaciones", 3] }, { $lte:
49 ["$numero_habitaciones", 5] }] }, then: "MEDIANO" }
50         ],
51         default: "GRANDE"
52       }
53     },
54     servicios_categoria: {
55       $switch: {
56         branches: [
57           { case: { $eq: [{ $size: "$servicios" }, 0] }, then: "SIN
58 SERVICIOS" },
59           { case: { $and: [{ $gte: [{ $size: "$servicios" }, 1] }, { $lte:
60 [{ $size: "$servicios" }, 10] }] }, then: "BASICO" },
61           { case: { $and: [{ $gte: [{ $size: "$servicios" }, 10] }, { $lte
62 : [{ $size: "$servicios" }, 25] }] }, then: "REGULAR" }
63         ],
64         default: "MUCHOS"
65       }
66     }
67   }
68 }

```

```
59     }
60   },
61   barrio_airbnb: 1,
62   nombre_airbnb: 1
63 }
64 },
65 {
66   $group: { // Se agrupan con respecto a las categorias creadas
67     _id: {
68       barrio: "$barrio_airbnb",
69       precio: "$precio_categoria",
70       puntuacion: "$puntuacion_categoria",
71       habitaciones: "$habitaciones_categoria",
72       servicios: "$servicios_categoria"
73     },
74     count: { $sum: 1 },
75     nombres_airbnb: { $addToSet: "$nombre_airbnb" } // Usamos $addToSet
76   } // para evitar nombres repetidos
77 },
78 {
79   $group: {
80     _id: "$_id.barrio",
81     categorias: {
82       $push: {
83         precio: "$_id.precio",
84         puntuacion: "$_id.puntuacion",
85         habitaciones: "$_id.habitaciones",
86         servicios: "$_id.servicios",
87         count: "$count",
88         nombres_airbnb: "$nombres_airbnb" // Incluimos el array de nombres
89       }
90     }
91   }
92 },
93 { $sort: { "_id": 1, "categorias.precio": 1, "categorias.puntuacion": 1, "categorias.habitaciones": 1, "categorias.servicios": 1 } }
94 ]
```



Figura 15: Categorización de alojamientos

3. Visualización de datos: Modelo de grafo

Se presenta un pequeño ejemplo del uso de un modelo de grafo en Neo4j, basado en los datos utilizados durante la tarea.

La estructura general de los nodos viene dada por:

- Local: Nombre, tipo de actividad, horario, licencia
- Terraza: Nombre, capacidad, estado de licencia
- Alojamiento: Nombre, precio, número de habitaciones, reseñas
- Barrio: Nombre

El modelo considera las siguientes relaciones:

- Ubicado_en: Conecta Local, Terraza y Alojamiento con su Barrio. Tiene el atributo *distancia*.
- Cercano_a: Conecta Local, Terraza y Alojamiento entre sí si están a menos de cierta distancia definida. Tiene el atributo *distancia*.
- Relacionado_con: conecta nodos de la misma categoría que comparten atributos comunes, como tipo de actividad o precio

3.1. Creación del grafo

Inicialmente creamos los nodos:

Nodos de barrios

```

1 CREATE (b1:Barrio {nombre: "Salamanca"})
2 CREATE (b2:Barrio {nombre: "Chamberi"})
3 CREATE (b3:Barrio {nombre: "Centro"})

```

Nodos de locales

```

1 CREATE (l1:Local {nombre: "Cafeteria Goya", tipo_actividad: "Restauracion",
   horario: "08:00-22:00", licencia: "Concedida"})
2 CREATE (l2:Local {nombre: "Libreria Central", tipo_actividad: "Cultura", horario:
   "10:00-20:00", licencia: "En tramite"})

```

Nodos de terrazas

```
1 CREATE (t1:Terraza {nombre: "Terraza Sol", capacidad: 30, estado_licencia: "
   Concedida"})
2 CREATE (t2:Terraza {nombre: "Terraza Norte", capacidad: 20, estado_licencia: "En
   tramite"})
```

Nodos de alojamiento

```
1 CREATE (a1:Alojamiento {nombre: "Airbnb Retiro", precio: 120, numero_habitaciones
   : 2, resenas: 45, servicios: ["WiFi", "Cocina"]})
2 CREATE (a2:Alojamiento {nombre: "Apartamento Goya", precio: 80,
   numero_habitaciones: 1, resenas: 0, servicios: ["Aire acondicionado"]})
```

Una vez creados los nodos establecemos las relaciones entre ellos:

Relacionamos locales, terrazas y alojamientos con barrios mediante “Ubicado_en”:

```
1 MATCH (b:Barrio {nombre: "Salamanca"}), (l:Local {nombre: "Cafeteria Goya"})
2 CREATE (l)-[:Ubicado_en {distancia: 0.5}]->(b)
3 MATCH (b:Barrio {nombre: "Centro"}), (a:Alojamiento {nombre: "Airbnb Retiro"})
4 CREATE (a)-[:Ubicado_en {distancia: 0.3}]->(b)
```

Ahora relacionamos entidades cercanas entre sí mediante “Cercano_a”:

```
1 MATCH (l:Local {nombre: "Cafeteria Goya"}), (t:Terraza {nombre: "Terraza Sol"})
2 CREATE (l)-[:Cercano_a {distancia: 0.2}]->(t)
3 MATCH (a:Alojamiento {nombre: "Airbnb Retiro"}), (l:Local {nombre: "Libreria
   Central"})
4 CREATE (a)-[:Cercano_a {distancia: 0.1}]->(l)
```

Finalmente relacionamos nodos que cuentan con atributos comunes:

```
1 MATCH (l:Local {tipo_actividad: "Restauracion"}), (t:Terraza {nombre: "Terraza
   Sol"})
2 CREATE (l)-[:Relacionado_con {motivo: "Ambiente de comida"}]->(t)
```

Como resultado obtenemos el grafo mostrado a continuación:

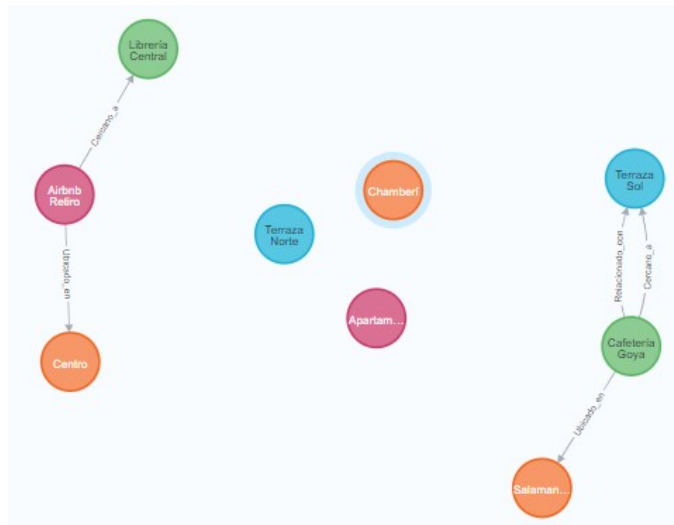


Figura 16: Grafo de comercio y alojamiento en Madrid

3.2. Consultas de validación

a - Locales y terrazas en Salamanca

```
1 MATCH (b:Barrio {nombre: "Salamanca"})<-[:Ubicado_en]-(n)
2 WHERE n:Local OR n:Terraza
3 RETURN n.nombre AS nombre, labels(n) AS tipo
```

neo4j\$ MATCH (b:Barrio {nombre: "Salamanca"})<-[:Ubicado_en]-(n) WHERE n:Local OR n:Terraza RETURN n.nombre AS nombre, labels(n) AS tipo

| nombre | tipo |
|------------------|-----------|
| "Cafetería Goya" | ["Local"] |

Figura 17: Locales y terrazas en Salamanca

b - Alojamientos cuyo precio supere 100 EUR

```
1 MATCH (a:Alojamiento)
2 WHERE a.precio > 100
3 RETURN a.nombre AS nombre, a.precio AS precio
```

neo4j\$ MATCH (a:Alojamiento) WHERE a.precio > 100 RETURN a.nombre AS nombre, a.precio AS precio

| nombre | precio |
|-----------------|--------|
| "Airbnb Retiro" | 120 |

Figura 18: Alojamientos cuyo precio supere 100 EUR

c - Barrios donde los alojamientos no cuentan con dormitorios

```
1 MATCH (a:Alojamiento)-[:Ubicado_en]->(b:Barrio)
2 WHERE a.numero_habitaciones = 0
3 RETURN b.nombre AS barrio, COUNT(a) AS total_alojamientos
```



Figura 19: Barrios donde los alojamientos no cuentan con dormitorios

d - Barrios donde los alojamientos no cuentan con reseñas

```
1 MATCH (a:Alojamiento)-[:Ubicado_en]->(b:Barrio)
2 WHERE a.resenas = 0
3 RETURN b.nombre AS barrio, COUNT(a) AS total_alojamientos
```

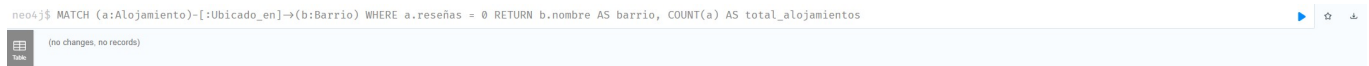


Figura 20: Barrios donde los alojamientos no cuentan con reseñas