Universidad de Córdoba

Máster en Inteligencia Computacional e Internet de las Cosas

MongoDB, Predicción de Actividades y Propiedades de Compuestos

Análisis, Diseño y Procesamiento de Datos Aplicados a las Ciencias y a las Tecnologías (ADP Marzo 2025)



Autor:

Marcos Rivera Gavilán

Profesor:

Gonzalo Cerrueda García

Índice general

Índice general				1
1.	Ejercicio 1			2
	1.1.	Objetiv	0	2
			sta de Estructura de Colecciones	2
			le Información de Prueba	4
			Colección departments	4
			Colección employees	4
	1.4.		as de Validación	6
			Listar todos los empleados	6
			Empleados con salario mayor a 2000	7
			Contar empleados por departamento	7
			Unir empleados con su departamento (\$lookup)	8
			Mostrar nombre del jefe de cada empleado	10
2.	Ejercicio 2			11
	•		0	11
		-	o a Realizar	
		Implementación en Python		

Capítulo 1

Ejercicio 1

1.1. Objetivo

El objetivo es afianzar los conocimientos sobre el uso de la base de datos NoSQL MongoDB, así como su modelado y acceso a la información.

1.2. Propuesta de Estructura de Colecciones

A partir del esquema E-R (Departamentos y Empleados, con la relación jefe-subordinado), se pueden diseñar dos colecciones principales en MongoDB:

- **departments**: almacena los datos de cada departamento (deptno, dname, loc, etc.).
- employees: almacena los datos de cada empleado (empno, ename, job, hiredate, sal, comm), con referencias a:
 - deptId: referencia al departamento correspondiente.
 - managerId: referencia al empleado que sea su jefe (opcional o null si es presidente).

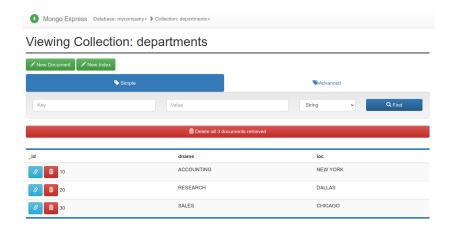


Figura 1.1: Colección departments en la base de datos.

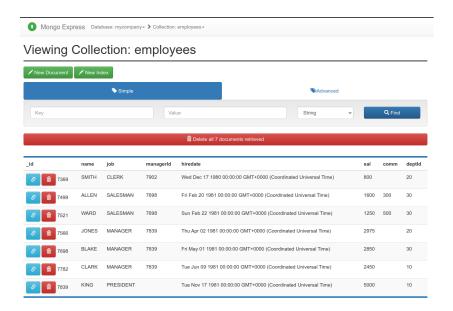


Figura 1.2: Colección employees en la base de datos.

Un ejemplo de documento en la colección departments podría ser:

Listing 1.1: Ejemplo de documento en departments

```
1 {
2    "_id": 10,
3    "dname": "ACCOUNTING",
4    "loc": "NEW YORK"
5 }
```

Un ejemplo de documento en la colección employees podría ser:

Listing 1.2: Ejemplo de documento en employees

```
1
2
      "_id": 7369,
      "name": "SMITH",
3
4
      "job": "CLERK",
5
      "managerId": 7902,
      "hiredate": ISODate("1980-12-17"),
6
      "sal": 800,
      "comm": null,
8
9
      "deptId": 20
10
```

1.3. Carga de Información de Prueba

Para contar con datos de validación, se pueden realizar inserciones en la base de datos mycompany de la siguiente forma:

1.3.1. Colección departments

Listing 1.3: Inserción en departments

1.3.2. Colección employees

Listing 1.4: Inserción en employees

```
db.employees.insertMany([
1
2
3
        _id: 7369,
        name: "SMITH",
4
5
        job: "CLERK",
        managerId: 7902,
6
7
        hiredate: new Date("1980-12-17"),
8
        sal: 800,
9
        comm: null,
10
        deptId: 20
     },
```

```
12
13
        _id: 7499,
        name: "ALLEN",
job: "SALESMAN",
14
15
        managerId: 7698,
        hiredate: new Date("1981-02-20"),
17
18
        sal: 1600,
19
        comm: 300,
        deptId: 30
20
21
22
23
        _id: 7521,
24
        name: "WARD",
25
        job: "SALESMAN",
^{26}
        managerId: 7698,
27
        hiredate: new Date("1981-02-22"),
28
        sal: 1250,
29
        comm: 500,
30
        deptId: 30
31
32
        _id: 7566,
33
        name: "JONES",
34
        job: "MANAGER",
35
36
        managerId: 7839,
        hiredate: new Date("1981-04-02"),
37
38
        sal: 2975,
39
        comm: null,
40
        deptId: 20
      },
41
42
        _id: 7698,
43
        name: "BLAKE",
44
        job: "MANAGER",
45
46
        managerId: 7839,
        hiredate: new Date("1981-05-01"),
47
48
        sal: 2850,
49
        comm: null,
        deptId: 30
50
51
52
        _id: 7782,
53
        name: "CLARK",
54
        job: "MANAGER",
55
56
        managerId: 7839,
        hiredate: new Date("1981-06-09"),
57
58
        sal: 2450,
59
        comm: null,
        deptId: 10
60
      },
61
62
      {
        _id: 7839,
63
64
        name: "KING",
        job: "PRESIDENT",
65
        // managerId: null o sin campo managerId
66
67
        hiredate: new Date("1981-11-17"),
68
        sal: 5000,
69
        comm: null,
70
        deptId: 10
      }
71
72
   ]);
```

1.4. Consultas de Validación

Con los datos cargados, se pueden realizar diversas consultas para verificar el acceso a la información:

1.4.1. Listar todos los empleados

Listing 1.5: Listado de empleados

```
db.employees.find().pretty();
2
3
    Γ
4
5
        _id: 7369,
        name: 'SMITH',
job: 'CLERK',
6
7
        managerId: 7902,
8
        hiredate: ISODate('1980-12-17T00:00:00.000Z'),
9
10
        sal: 800,
        comm: null,
11
12
        deptId: 20
13
14
15
        _id: 7499,
        name: 'ALLEN',
16
        job: 'SALESMAN',
17
18
        managerId: 7698,
19
        hiredate: ISODate('1981-02-20T00:00:00.000Z'),
        sal: 1600,
20
21
        comm: 300,
22
        deptId: 30
23
24
25
        _id: 7521,
26
        name: 'WARD',
        job: 'SALESMAN',
27
28
        managerId: 7698,
29
        hiredate: ISODate('1981-02-22T00:00:00.000Z'),
        sal: 1250,
30
31
        comm: 500,
32
        deptId: 30
      },
33
34
        _id: 7566,
35
        name: 'JONES'
36
37
        job: 'MANAGER',
        managerId: 7839,
38
        hiredate: ISODate('1981-04-02T00:00:00.000Z'),
39
40
        sal: 2975,
41
        comm: null,
42
        deptId: 20
43
44
45
        _id: 7698,
        name: 'BLAKE',
46
        job: 'MANAGER'
47
48
        managerId: 7839,
        hiredate: ISODate('1981-05-01T00:00:00.000Z'),
49
50
        sal: 2850,
        comm: null,
```

```
deptId: 30
52
53
      },
54
        _id: 7782,
55
        name: 'CLARK',
57
        job: 'MANAGER',
58
        managerId: 7839,
        hiredate: ISODate('1981-06-09T00:00:00.000Z'),
59
60
        sal: 2450,
61
        comm: null,
        deptId: 10
62
63
64
65
        _id: 7839,
66
        name: 'KING',
67
        job: 'PRESIDENT',
        hiredate: ISODate('1981-11-17T00:00:00.000Z'),
68
69
        sal: 5000,
70
        comm: null,
71
        deptId: 10
72
73
   ]
```

1.4.2. Empleados con salario mayor a 2000

Listing 1.6: Empleados con sal ¿2000

1.4.3. Contar empleados por departamento

Listing 1.7: Total de empleados por deptId

1.4.4. Unir empleados con su departamento (\$lookup)

Listing 1.8: Unión empleado-departamento

```
db.employees.aggregate([
1
2
3
        $lookup: {
4
          from: "departments",
          localField: "deptId",
5
          foreignField: "_id",
6
          as: "deptInfo"
7
8
9
     }
10
   ]);
11
12
   [
13
        _id: 7369,
14
        name: 'SMITH',
15
        job: 'CLERK',
16
        managerId: 7902,
17
        hiredate: ISODate('1980-12-17T00:00:00.000Z'),
18
19
        sal: 800,
20
        comm: null,
21
        deptId: 20,
22
        deptInfo: [ { _id: 20, dname: 'RESEARCH', loc: 'DALLAS' } ]
23
24
25
        _id: 7499,
        name: 'ALLEN',
26
        job: 'SALESMAN',
27
28
        managerId: 7698,
        hiredate: ISODate('1981-02-20T00:00:00.000Z'),
29
30
        sal: 1600,
        comm: 300,
31
32
        deptId: 30,
33
        deptInfo: [ { _id: 30, dname: 'SALES', loc: 'CHICAGO' } ]
34
35
        _id: 7521,
36
        name: 'WARD',
37
38
        job: 'SALESMAN',
39
        managerId: 7698,
40
        hiredate: ISODate('1981-02-22T00:00:00.000Z'),
41
        sal: 1250,
42
        comm: 500,
43
        deptId: 30,
        deptInfo: [ { _id: 30, dname: 'SALES', loc: 'CHICAGO' } ]
44
45
     },
46
        _id: 7566,
47
48
        name: 'JONES',
49
        job: 'MANAGER'
        managerId: 7839,
50
51
        hiredate: ISODate('1981-04-02T00:00:00.000Z'),
52
        sal: 2975,
        comm: null,
53
54
        deptId: 20,
55
        deptInfo: [ { _id: 20, dname: 'RESEARCH', loc: 'DALLAS' } ]
56
57
58
        _id: 7698,
```

```
name: 'BLAKE',
59
        job: 'MANAGER',
60
        managerId: 7839,
hiredate: ISODate('1981-05-01T00:00:00.000Z'),
61
62
63
        sal: 2850,
64
        comm: null,
65
        deptId: 30,
66
        deptInfo: [ { _id: 30, dname: 'SALES', loc: 'CHICAGO' } ]
67
      },
68
        _id: 7782, name: 'CLARK',
69
70
        job: 'MANAGER',
71
        managerId: 7839,
72
        hiredate: ISODate('1981-06-09T00:00:00.000Z'),
73
74
        sal: 2450,
        comm: null,
75
76
        deptId: 10,
77
        deptInfo: [ { _id: 10, dname: 'ACCOUNTING', loc: 'NEW YORK' } ]
78
79
        _id: 7839,
80
        name: 'KING',
81
        job: 'PRESIDENT',
82
        hiredate: ISODate('1981-11-17T00:00:00.000Z'),
83
84
        sal: 5000,
85
        comm: null,
86
        deptId: 10,
87
        deptInfo: [ { _id: 10, dname: 'ACCOUNTING', loc: 'NEW YORK' } ]
      }
88
89
   ]
```

1.4.5. Mostrar nombre del jefe de cada empleado

Listing 1.9: Relación jefe-subordinado

```
db.employees.aggregate([
 2
 3
                $lookup: {
                    from: "employees",
localField: "managerId",
 4
 5
  6
                    foreignField: "_id",
 7
                    as: "managerInfo"
  8
 9
10
11
                $project: {
12
                    name: 1,
13
                    job: 1,
                    manager: { $arrayElemAt: ["$managerInfo.name", 0] }
14
15
           }
16
17
       ]);
18
19
           { _id: 7369, name: 'SMITH', job: 'CLERK' },
{ _id: 7499, name: 'ALLEN', job: 'SALESMAN', manager: 'BLAKE' },
20
21
           { _id: 7499, name: 'ALLEN', job: 'SALESMAN', manager: 'BLAKE' }, { _id: 7521, name: 'WARD', job: 'SALESMAN', manager: 'BLAKE' }, { _id: 7566, name: 'JONES', job: 'MANAGER', manager: 'KING' }, { _id: 7698, name: 'BLAKE', job: 'MANAGER', manager: 'KING' }, { _id: 7782, name: 'CLARK', job: 'MANAGER', manager: 'KING' }, { _id: 7839, name: 'KING', job: 'PRESIDENT' }
23
24
25
26
27
       ]
```

Capítulo 2

Ejercicio 2

2.1. Objetivo

El objetivo de esta práctica es afianzar los conocimientos impartidos en la parte teórica sobre la caracterización y predicción de actividades o propiedades de compuestos. En concreto, se pretende realizar la predicción (clasificación) del tipo de actividad biológica (Activo "1"; Inactivo "0") a partir del fingerprint molecular de cada compuesto.

2.2. Trabajo a Realizar

- 1. Conectar a las bases de datos CDS16 y CDS29 (u otra base de datos similar) que contienen las colecciones mfp_counts y molecules.
- 2. Extraer los datos de la colección molecules, donde el campo class indica la actividad biológica.
- 3. Reconstruir la matriz de variables predictoras (fingerprint de 1024 bits, en columnas FP1, FP2, ..., FP1024) a partir de las posiciones que están a "1".
- 4. Eliminar columnas en las que todos los valores sean "0".
- 5. Utilizar alguno de los modelos de scikit-learn (RandomForestClassifier, SVC, etc.) para entrenar y predecir la actividad.
- Evaluar el desempeño con métricas como Accuracy, Kappa, AUC, GMean, etc.



Figura 2.1: Base de Datos CDS16



Figura 2.2: Base de Datos CDS29

2.3. Implementación en Python

Listing 2.1: Ejemplo unificado para la conexión y predicción

```
1 import pymongo
   import numpy as np
3 import pandas as pd
4 import matplotlib.pyplot as plt
   from sklearn.model_selection import train_test_split, GridSearchCV,
         StratifiedKFold
    {\tt from \ sklearn.ensemble \ import \ RandomForestClassifier}
7
8
   from sklearn.metrics import accuracy_score, cohen_kappa_score,
        roc_auc_score
9
10
   # Helper function to load and process a dataset from a given
        collection.
12
13 def load_dataset(collection, num_bits=1024):
        docs = list(collection.find({}))
14
        df_rows = []
15
16
17
        for doc in docs:
            # Create a fingerprint array of 1024 bits initialized to 0.
18
19
            fp_row = np.zeros(num_bits, dtype=int)
20
21
            # Update the array if the document contains a valid
                fingerprint list.
            if "fingerprint" in doc and doc["fingerprint"]:
22
                for pos in doc["fingerprint"]:
23
                     if 0 <= pos < num_bits:
24
25
                         fp_row[pos] = 1
26
            \mbox{\tt\#} Build a dictionary with "class" and each fingerprint bit.
27
28
            row_data = {"class": doc.get("class", 0)} # default to 0
                if "class" is missing
29
            for i in range(num_bits):
30
                row_data[f"FP{i+1}"] = fp_row[i]
31
            df_rows.append(row_data)
32
33
        # Convert list of dictionaries to a DataFrame.
34
        df = pd.DataFrame(df_rows)
35
        print("Before filtering, DataFrame shape:", df.shape)
36
37
        # Filter out fingerprint columns that are all zeros.
38
        filtered_df = df.loc[:, (df != 0).any(axis=0)]
39
        \mbox{\tt\#} If filtering removes all fingerprint columns (leaving only \mbox{\tt\#}
            class"), revert.
40
        if len(filtered_df.columns) <= 1:</pre>
            print("Warning: No non-zero fingerprint columns found;
41
                reverting to original DataFrame.")
42
            filtered_df = df
43
        print("After filtering, DataFrame shape:", filtered_df.shape)
45
        return filtered_df
46
```

```
47 #
              ______
48 # Helper function to perform grid search for hyperparameter tuning.
50
    def grid_search_rf(X_train, y_train):
51
        param_grid = {
            'n_estimators': [100, 200],
52
            'max_depth': [None, 10, 20],
'min_samples_split': [2, 5],
53
54
55
        rf = RandomForestClassifier(random_state=42, class_weight=')
56
           balanced')
57
        cv = StratifiedKFold(n_splits=5, shuffle=True, random_state=42)
        grid_search = GridSearchCV(rf, param_grid, cv=cv, scoring='
58
            roc_auc', n_jobs=-1)
        grid_search.fit(X_train, y_train)
print("Best parameters:", grid_search.best_params_)
59
60
61
        return grid_search.best_estimator_
62
63 #
64 # Helper function to train and evaluate a model on a DataFrame.
65 #
   def train_and_evaluate(df):
67
        \# Separate features (X) and target (y)
        X = df.drop('class', axis=1)
68
        y = df['class']
69
70
71
        # Print class distribution to assess imbalance.
        print("Class distribution:")
72
73
        print(y.value_counts())
74
75
        # Split into training and testing sets (using stratification).
76
        X_train, X_test, y_train, y_test = train_test_split(
            X, y, test_size=0.2, random_state=42, stratify=y
77
78
        )
79
        \# Hyperparameter tuning with GridSearchCV.
80
81
        best_rf = grid_search_rf(X_train, y_train)
82
83
        # Train the tuned classifier on the full training set.
84
        best_rf.fit(X_train, y_train)
85
86
        # Predict on test data.
87
        y_pred = best_rf.predict(X_test)
        y_pred_prob = best_rf.predict_proba(X_test)[:, 1]
88
89
90
        # Calculate evaluation metrics.
91
        acc = accuracy_score(y_test, y_pred)
92
        kappa = cohen_kappa_score(y_test, y_pred)
93
        auc = roc_auc_score(y_test, y_pred_prob)
94
95
        \mbox{\tt\#} Plot histogram of predicted probabilities for diagnostics.
96
        plt.hist(y_pred_prob, bins=20)
97
        plt.xlabel("Predicted Probability")
        plt.ylabel("Frequency")
```

```
99
        plt.title("Distribution of Predicted Probabilities")
100
        plt.show()
101
        return best_rf, acc, kappa, auc
102
103
104 #
105 # 1) Connect to MongoDB.
106
107
    client = pymongo.MongoClient("mongodb://admin:1234@localhost:27017/
108
109
110 # 2) Select the CDS16 database and its collections.
111 #
112 db_cds16 = client["CDS16"]
   collection_molecules_16 = db_cds16["molecules"]
113
    collection_mfp_counts_16 = db_cds16["mfp_counts"]
114
115
116 print(f"CDS16 - molecules: {collection_molecules_16.count_documents
       ({})} documentos.")
    print(f"CDS16 - mfp_counts: {collection_mfp_counts_16.
        count_documents({})} documentos.")
118
119 #
120 # 3) Select the CDS29 database and its collections.
121
        ______
    db_cds29 = client["CDS29"]
122
    collection_molecules_29 = db_cds29["molecules"]
    collection_mfp_counts_29 = db_cds29["mfp_counts"]
124
125
126
    print(f"CDS29 - molecules: {collection_molecules_29.count_documents
        ({})} documentos.")
127
    print(f"CDS29 - mfp_counts: {collection_mfp_counts_29.
        count_documents({})} documentos.")
128
129
130
   # 4) Load and process datasets from both CDS16 and CDS29.
131
132 df_cds16 = load_dataset(collection_molecules_16)
    df_cds29 = load_dataset(collection_molecules_29)
134
135
   # Ensure that the feature matrix is not empty.
136 if df_cds16.drop('class', axis=1).empty:
        raise ValueError("Feature matrix for CDS16 is empty. Please
137
            check your data ingestion and filtering steps.")
138 if df_cds29.drop('class', axis=1).empty:
```

```
139
           raise ValueError("Feature matrix for CDS29 is empty. Please
              check your data ingestion and filtering steps.")
140
141 #
142\, # 5) Train and evaluate the model on both datasets.
143 #
144 rf_cds16, acc_cds16, kappa_cds16, auc_cds16 = train_and_evaluate(
           df_cds16)
     rf_cds29, acc_cds29, kappa_cds29, auc_cds29 = train_and_evaluate(
           df_cds29)
146
147 #
148\, # 6) Print and compare performance metrics.
149 #
           ______
150 comparison_df = pd.DataFrame({
151 "Dataset": ["CDS16", "CDS29"],
           "Accuracy": [acc_cds16, acc_cds29],
"Kappa": [kappa_cds16, kappa_cds29],
152
153
           "AUC": [auc_cds16, auc_cds29]
154
155 })
156
157 print("\nComparison of model performance:")
158 \quad \mathtt{print} (\mathtt{comparison\_df})
159
160 Comparison of model performance:

        161
        Dataset
        Accuracy
        Kappa
        AUC

        162
        0
        CDS16
        0.470588
        0.0
        0.5

        163
        1
        CDS29
        0.803991
        0.0
        0.5
```