



UNIVERSIDAD DE CÓRDOBA

Departamento de Informática y Análisis Numérico

Ingeniería del Software, Conocimiento y Bases de Datos

MASTERES UNIVERSITARIOS-UNIVERSIDAD DE CÓRDOBA

ASIGNATURA:

Análisis, Diseño y Procesamiento de Datos Aplicados a las Ciencias y a las Tecnologías

Nombre: Lebbihi

Apellido: Mhammed

Profesores: 1- Gonzalo Cerreuil García

2- Domingo Ortiz Boyer

3- Juan Antonio Romero

PRÁCTICAS (APARTADO III)

Etapla 1: Diseño de la Estructura de Datos

Objetivo: Definir cómo se almacenarán y organizarán los datos dentro de MongoDB.

Acciones:

Crear una colección departments para los departamentos, con campos para el número del departamento, nombre y localización.

Crear una colección employees para los empleados, con campos para el número de empleado, nombre, puesto, fecha de contratación, salario, comisión y el número del departamento.

Considerar las relaciones entre los departamentos y los empleados para determinar si se necesitan referencias o documentos embebidos.

Etapla 2: Implementación de la Estructura de Datos

Objetivo: Establecer las colecciones y sus relaciones en la base de datos MongoDB.

Acciones:

Usar comandos de MongoDB para crear las colecciones departments y employees.

Definir los índices necesarios para optimizar las consultas, especialmente para campos que se usarán frecuentemente en búsquedas y relaciones.

Etapla 3: Inserción de Datos

Objetivo: Poblar las colecciones con datos de prueba que sean representativos de los datos reales.

Acciones:

Generar o insertar datos de prueba en la colección departments.

Generar o insertar datos de prueba en la colección employees, asegurándose de que los datos reflejen las relaciones adecuadas con la colección departments.

Etapas 5: Optimización y Pruebas

Objetivo: Asegurarse de que la solución sea eficiente y se ajuste a las necesidades de la aplicación.

Acciones:

Analizar el rendimiento de las consultas y realizar ajustes en los índices si es necesario.

Revisar el esquema de datos y realizar cualquier ajuste para mejorar la estructura basándose en los patrones de acceso a los datos.

Estas etapas cubren el proceso de diseño e implementación de una solución de base de datos usando MongoDB para la parte práctica de un curso, con un enfoque en el aprendizaje teórico y la aplicación práctica de estos conocimientos.

File Edit Selection View Go Run ... Aportado

MONGODB

- CONNECTIONS
 - clusters1-301rpx.mongodb.net
 - cluster31-301rpx.mongodb.net
 - localhost:27017 connected
 - admin
 - Aportado
 - departments
 - employees
 - config
- PLAYGROUNDS

No MongoDB playground files found in the workspace.

Create New Playground
- HELP AND FEEDBACK
 - What's New
 - Extension Documentation
 - MongoDB Documentation
 - Suggest a Feature
 - Report a Bug
 - Create Free Atlas Cluster

Aportado.py

```
1 from pymongo import MongoClient
2 from random import randint
3 import datetime
4
5 # Connexion à MongoDB
6 client = MongoClient('mongodb://localhost:27017/')
7 db = client['Aportado']
8
9 # Collection Departments
10 departments = db['departments']
11 # Collection Employees
12 employees = db['employees']
13
14 # Créer 10 départements
15 for i in range(1, 11):
16     department = {
17         "DEPTNO": str(i),
18         "DNAME": f"Departement_{i}",
19         "LOC": f"Emplacement_{i}"
20     }
21     departments.insert_one(department)
22
23 # Créer 10 employés et les assigner à des départements aléatoires
24 for j in range(1, 11):
25     employee = {
26         "EMPNO": str(1000 + j),
27         "ENAME": f"Employe_{j}",
28         "JOB": f"Poste_{j}",
29         "HIREDATE": datetime.datetime.now(),
30         "SAL": randint(3000, 6000)
31     }
```

MongoDB connection successful.

Ln 59, Col 1 Spaces: 4 UTF-8 CRLF Python 3.11.4 64-bit 22:47 16/03/2024

File Edit Selection View Go Run ... Aportado

MONGODB

- CONNECTIONS
 - clusters1-301rpx.mongodb.net
 - cluster31-301rpx.mongodb.net
 - localhost:27017 connected
 - admin
 - Aportado
 - departments
 - employees
 - config
- PLAYGROUNDS

No MongoDB playground files found in the workspace.

Create New Playground
- HELP AND FEEDBACK
 - What's New
 - Extension Documentation
 - MongoDB Documentation
 - Suggest a Feature
 - Report a Bug
 - Create Free Atlas Cluster

Aportado.py

```
37 employees.insert_one(employee)
38
39 # Trouver tous les employés d'un département
40 for employee in employees.find({"DEPTNO": "10"}):
41     print(employee)
42
43 # Trouver le chef d'un employé donné
44 boss = employees.find_one({"EMPNO": "7934"})['MANAGER']
45 print(employees.find_one({"EMPNO": boss}))
46
47 # Agréger les salaires par département
48 pipeline = [
49     {
50         "$group": {
51             "_id": "$DEPTNO",
52             "total_salary": {"$sum": "$SAL"}
53         }
54     }
55 ]
```

PROBLEMS OUTPUT DEBUG CONSOLE TERMINAL PORTS Python + - Python 3.11.4 64-bit

DATE: datetime.datetime(2024, 1, 15, 22, 16, 20, 45000), 'SAL': 58122, 'DEPTNO': '10', 'MANAGER': '1002'}
{'_id': ObjectId('65f6056d3e88639d0359186d'), 'EMPNO': '7839', 'ENAME': 'King', 'JOB': 'President', 'HIREDAT
{'_id': '6', 'total_salary': 140416}
{'_id': '8', 'total_salary': 42684}
{'_id': '4', 'total_salary': 131725}
{'_id': '10', 'total_salary': 93255}
{'_id': '2', 'total_salary': 35678}
{'_id': '7', 'total_salary': 133157}
{'_id': '5', 'total_salary': 66769}
{'_id': '1', 'total_salary': 41170}
{'_id': '3', 'total_salary': 140303}
{'_id': '9', 'total_salary': 52797}

Ln 59, Col 1 Spaces: 4 UTF-8 CRLF Python 3.11.4 64-bit 22:48 16/03/2024

MongoDB Compass - localhost:27017/Aportado

Connect Edit View Help

localhost:27017

My Queries Performance Databases Search

Aportado

- departments
- employees
- Milan_CDR_db
 - Milan_CDR_o
- admin
- config
- local

+ Create collection Refresh

View Sort by Collection Name

departments

Storage size:	Documents:	Avg. document size:	Indexes:	Total index size:
20.48 kB	21	84.00 B	1	36.86 kB

employees

Storage size:	Documents:	Avg. document size:	Indexes:	Total index size:
20.48 kB	22	133.00 B	1	36.86 kB

MONGOSH

Buscar

16°C 22:49 16/03/2024

MongoDB Compass - localhost:27017/Aportado.departments

Connect Edit View Collection Help

localhost:27017

My Queries Performance Databases Search

Aportado

- departments
- employees
- Milan_CDR_db
 - Milan_CDR_o
- admin
- config
- local

Aportado.departments 21 1 DOCUMENTS INDEXES

Documents Aggregations Schema Indexes Validation

Filter Type a query: { field: 'value' } or Generate query

EXPLAIN Reset Find Options

ADD DATA EXPORT DATA UPDATE DELETE

1 - 20 of 21

```
{ "_id": ObjectId('65f607f49162e7a9b9aeaaa5'), DEPTNO: "1", DNAME: "Departement_1", LOC: "Emplacement_1" }
```

```
{ "_id": ObjectId('65f607f49162e7a9b9aeaaa6'), DEPTNO: "2", DNAME: "Departement_2", LOC: "Emplacement_2" }
```

```
{ "_id": ObjectId('65f607f49162e7a9b9aeaaa7'), DEPTNO: "3", DNAME: "Departement_3", LOC: "Emplacement_3" }
```

MONGOSH

2: APLICACIONES CIENTÍFICAS Y EMPRESARIALES (I)

Etapas 1: Recopilación de Datos

Extraer los datos de compuestos de una base de datos, donde cada compuesto tiene un conjunto de características (fingerprint molecular) y una etiqueta que indica su actividad biológica (activa o inactiva).

Etapas 2: Limpieza y Preparación de Datos

Limpiar los datos para la modelización. Esto podría incluir la eliminación de columnas que sólo contienen ceros y la transformación de datos categóricos en numéricos si es necesario.

Etapas 3: División de Datos

Dividir el conjunto de datos en partes de entrenamiento y prueba para poder evaluar el modelo de manera justa.

Etapas 4: Desarrollo del Modelo

Desarrollar un modelo de clasificación que utilice las características de los compuestos para predecir su actividad biológica.

Etapas 5: Entrenamiento del Modelo

Entrenar el modelo con el conjunto de datos de entrenamiento, ajustando los parámetros según sea necesario para mejorar el rendimiento.

Etapas 6 y 7: Pruebas y Validación

Probar el modelo con el conjunto de datos de prueba para evaluar su capacidad de generalización y su rendimiento con datos no vistos durante el entrenamiento.

Etapas 8 y 9: Evaluación de Rendimiento

Evaluar el rendimiento del modelo utilizando métricas como la precisión, la sensibilidad, la especificidad y el área bajo la curva ROC (AUC-ROC).

Etapas 10 y 11: Ajuste Fino

Realizar ajustes finos en el modelo según sea necesario para mejorar el rendimiento, lo cual puede incluir la selección de características, la ingeniería de características o la optimización de hiperparámetros.

Etapas 12 y 13: Interpretación y Conclusión

Interpretar los resultados del modelo, entender la importancia de las diferentes características y sacar conclusiones que puedan aplicarse en contextos científicos o empresariales.

Etapas 14 y 15: Documentación y Reporte

Documentar el proceso y los hallazgos y preparar un informe detallado para los interesados o para la publicación de los resultados de la investigación.

Este resumen cubre el proceso general de desarrollo y evaluación de un modelo de clasificación en el contexto de aplicaciones científicas y empresariales, desde la recopilación de datos hasta la documentación final y la generación de informes.

The screenshot shows the MongoDB Playground interface. On the left, the 'CONNECTIONS' panel lists several clusters, with 'localhost27017' selected and connected. Below it, the 'PLAYGROUNDS' panel indicates no files are found in the workspace. The main editor area, titled 'from sklearn.py', contains a Python script for phase 2. The script imports necessary libraries, simulates data loading from MongoDB, creates a synthetic DataFrame with 100 samples and 1024 features, and splits the data into training and testing sets using `train_test_split`. The status bar at the bottom shows the file is at line 5, column 13, with 4 spaces, in UTF-8 encoding, using Python 3.11.4 64-bit.

```
from sklearn.py > ...
phase 2 > from sklearn.py > ...
1 from sklearn.model_selection import train_test_split
2 from sklearn.ensemble import RandomForestClassifier
3 from sklearn.metrics import accuracy_score, roc_auc_score
4 import numpy as np
5 import pandas as pd
6
7 # Simular la carga de datos. En tu caso, cargarías los datos de MongoDB.
8 # Por ahora, voy a crear un DataFrame de ejemplo con datos ficticios.
9 # Tu DataFrame debe ser creado basándote en tus datos reales de MongoDB.
10
11 # Crear un DataFrame de ejemplo
12 num_samples = 100 # Supongamos que tienes 100 compuestos
13 num_features = 1024 # Supongamos 1024 fingerprints
14
15 # Generar datos aleatorios
16 X = np.random.randint(2, size=(num_samples, num_features))
17 y = np.random.randint(2, size=(num_samples))
18
19 # Convertir a DataFrame y asignar clase
20 df = pd.DataFrame(X, columns=[f'FP{i+1}' for i in range(num_features)])
21 df['class'] = y
22
23 # Eliminar columnas donde todos los valores sean '0'
24 df = df.loc[:, (df != 0).any(axis=0)]
25
26 # Dividir los datos en entrenamiento y prueba
27 X = df.drop('class', axis=1)
28 y = df['class']
29 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3, random_state=42)
30
```

This screenshot shows the continuation of the Python script in the MongoDB Playground. The script proceeds with training a `RandomForestClassifier` model on the training data, making predictions on the test set, and calculating performance metrics including accuracy and ROC AUC. The status bar at the bottom indicates the script is now at line 5, column 13, with 4 spaces, in UTF-8 encoding, using Python 3.11.4 64-bit.

```
from sklearn.py > ...
phase 2 > from sklearn.py > ...
17 y = np.random.randint(2, size=(num_samples))
18
19 # Convertir a DataFrame y asignar clase
20 df = pd.DataFrame(X, columns=[f'FP{i+1}' for i in range(num_features)])
21 df['class'] = y
22
23 # Eliminar columnas donde todos los valores sean '0'
24 df = df.loc[:, (df != 0).any(axis=0)]
25
26 # Dividir los datos en entrenamiento y prueba
27 X = df.drop('class', axis=1)
28 y = df['class']
29 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3, random_state=42)
30
31 # Crear y entrenar el modelo de clasificación
32 clf = RandomForestClassifier(n_estimators=100, random_state=42)
33 clf.fit(X_train, y_train)
34
35 # Realizar predicciones en el conjunto de prueba
36 y_pred = clf.predict(X_test)
37
38 # Calcular e imprimir métricas de rendimiento
39 accuracy = accuracy_score(y_test, y_pred)
40 roc_auc = roc_auc_score(y_test, clf.predict_proba(X_test)[:, 1])
41 print(f'Accuracy: {accuracy:.2f}')
42 print(f'ROC AUC: {roc_auc:.2f}')
43
```


Visual Studio Code interface showing a Python script for data analysis using sklearn and pandas. The script is named "from sklearn.py" and is located in the "Aportado" folder.

Left Panel (MongoDB Explorer):

- CONNECTIONS
 - clusters1301rpn.mongodb.net
 - cluster313oi9tph.mongodb.net
 - localhost27017 connected
 - admin
 - Aportado
 - departments
 - employees
 - config
- PLAYGROUNDS
 - No MongoDB playground files found in the workspace.
 - Create New Playground
- HELP AND FEEDBACK
 - What's New
 - Extension Documentation
 - MongoDB Documentation
 - Suggest a Feature
 - Report a Bug
 - Create Free Atlas Cluster

Main Editor (from sklearn.py):

```
phase 2 > from sklearn.py > ...
17 y = np.random.randint(2, size=(num_samples))
18
19 # Convertir a DataFrame y asignar clase
20 df = pd.DataFrame(X, columns=[f'FP{i+1}' for i in range(num_features)])
21 df['class'] = y
22
23 # Eliminar columnas donde todos los valores sean '0'
24 df = df.loc[:, (df != 0).any(axis=0)]
25
26 # Dividir los datos en entrenamiento y prueba
27 X = df.drop('class', axis=1)
28 y = df['class']
29 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3, random_state=42)
30
31 # Crear y entrenar el modelo de clasificación
32 clf = RandomForestClassifier(n_estimators=100, random_state=42)
33 clf.fit(X_train, y_train)
```

Terminal:

```
> & C:/Users/Windows10/AppData/Local/Programs/Python/Python311/python.exe "c:/Users/Windows10/Desktop/Master/ADP/Aportado/phase 2/from sklearn.py"
Accuracy: 0.47
ROC AUC: 0.42
PS C:\Users\Windows10\Desktop\Master\ADP\Aportado>
```

Bottom Panel:

- PROBLEMS
- OUTPUT
- DEBUG CONSOLE
- TERMINAL
- PORTS
 - powershell
 - Python
 - powershell

Status Bar: Ln 5, Col 13 Spaces: 4 UTF-8 Python 3.11.4 64-bit 23:17 16/03/2024