

## Mongodb Handler Class

```
1 | # Author : Ashantha Rosary
2 | from pymongo import MongoClient
3 | from pyspark.sql import SparkSession
4 | from pyspark.sql.types import StructType, StructField, StringType, IntegerType
5 | import pprint
6 |
7 | class MongoDBHandler:
8 |     """
9 |     A class to handle operations with MongoDB.
10 |    """
11 |    def __init__(self, uri, database, collection):
12 |        """
13 |        Initializes the MongoDB client with connection details.
14 |
15 |        Parameters:
16 |        - uri: str, URI for the MongoDB database.
17 |        - database: str, name of the MongoDB database.
18 |        - collection: str, name of the MongoDB collection.
19 |        """
20 |        self.uri = uri
21 |        self.database = database
22 |        self.collection = collection
23 |        self.client = MongoClient(self.uri)
24 |        self.db = self.client[self.database]
25 |        self.collection = self.db[self.collection]
26 |
27 |    def list_documents(self, limit=3):
28 |        """
29 |        Lists a specified number of documents from the MongoDB collection.
30 |
31 |        Parameters:
32 |        - limit: int, number of documents to list.
33 |        """
34 |        print(f"Listing {limit} documents in the {self.collection.name} collection: ")
35 |        head_review = self.collection.find().limit(limit)
36 |        pprint.pprint(list(head_review))
37 |
38 |    def retrieve_data(self):
39 |        """
40 |        Retrieves all data from the MongoDB collection.
41 |
42 |        Returns:
43 |        - list of records retrieved from MongoDB.
44 |        """
45 |        mongo_data = list(self.collection.find())
46 |        return mongo_data
47 |
```

```

48 def convert_to_dataframe(self, mongo_data, spark_session, schema):
49     """
50     Converts MongoDB records into a PySpark DataFrame.
51
52     Parameters:
53     - mongo_data: list of records retrieved from MongoDB.
54     - spark_session: SparkSession object to create DataFrame.
55     - schema: StructType, schema definition for the DataFrame.
56
57     Returns:
58     - PySpark DataFrame.
59     """
60     # Remove the _id field from each record
61     for record in mongo_data:
62         if '_id' in record:
63             del record['_id']
64
65     df = spark_session.createDataFrame(mongo_data, schema=schema)
66     return df
67
68 def close(self):
69     """
70     Closes the MongoDB client connection.
71     """
72     if self.client is not None:
73         self.client.close()

```