

Hbase Handler Class

```
1  ## Author: Goh Boon Xiang
2  import happybase
3  import uuid
4  from pyspark.sql import SparkSession, Row
5  from pyspark.sql.types import StructType, StructField, StringType, IntegerType,
   ArrayType, DoubleType
6  from pyspark.ml.linalg import Vectors, VectorUDT
7
8  class HBaseHandler:
9      """
10     A class to handle operations with HBase.
11     """
12     def __init__(self, host='localhost', port=9090):
13         """
14         Initializes the HBase connection.
15
16         Parameters:
17         - host: str, HBase server host.
18         - port: int, HBase server port.
19         """
20         self.host = host
21         self.port = port
22         self.connection = happybase.Connection(self.host, port=self.port)
23         self.connection.open()
24
25     def list_tables(self):
26         """
27         Lists all the tables in HBase.
28         """
29         print("Available tables:", self.connection.tables())
```

```

31 def delete_table(self, table_name):
32     """
33     Deletes a specified table from HBase.
34
35     Parameters:
36     - table_name: str, the name of the table to delete.
37     """
38     if table_name.encode('utf-8') in self.connection.tables():
39         try:
40             self.connection.disable_table(table_name)
41             print(f"Table '{table_name}' disabled successfully.")
42         except happybase.hbase.ttypes.IOException as e:
43             if 'TableNotDisabledException' in str(e):
44                 print(f"Table '{table_name}' is already disabled.")
45             else:
46                 raise e
47
48         self.connection.delete_table(table_name)
49         print(f"Table '{table_name}' deleted successfully.")
50     else:
51         print(f"Table '{table_name}' does not exist.")
52
53     self.list_tables()
54
55 def create_table(self, table_name, families):
56     """
57     Creates a table in HBase.
58
59     Parameters:
60     - table_name: str, the name of the table to create.
61     - families: dict, the column families to be created.
62     """
63     self.connection.create_table(table_name, families)
64     print(f"Table '{table_name}' created successfully.")
65     self.list_tables()
66
67 def generate_row_key(self, row):
68     """
69     Generates a unique row key for storing in HBase.
70
71     Parameters:
72     - row: Row object containing the data.
73
74     Returns:
75     - str, the generated row key.
76     """
77     unique_id = uuid.uuid4()
78     key_hash = hash((row['Sentiment'], row['Review']))
79     return f"{unique_id}_{key_hash}"

```

```

81     def save_to_hbase(self, df, table_name):
82         """
83         Saves data from a DataFrame to HBase.
84
85         Parameters:
86         - df: DataFrame, the data to be stored.
87         - table_name: str, the name of the HBase table.
88         """
89         table = self.connection.table(table_name)
90         total = 0
91
92         for row in df.collect():
93             row_key = self.generate_row_key(row)
94             data_to_store = {
95                 b'cf1:Review': row['Review'].encode() if row['Review'] else b'',
96                 b'cf2:Sentiment': str(row['Sentiment']).encode() if row['Sentiment'] is
100 not None else b'',
97                 b'cf3:SkuInfo_index': str(row['SkuInfo_index']).encode() if
101 row['SkuInfo_index'] is not None else b'',
98                 b'cf4:tokens': ','.join(row['tokens']).encode() if row['tokens'] else
102 b'',
99                 b'cf5:number_of_tokens': str(row['number_of_tokens']).encode() if
103 row['number_of_tokens'] is not None else b'',
104
105             }
106             table.put(row_key.encode(), data_to_store)
107             total += 1
108
109         print(f'Data successfully stored in HBase with {total} records')

```

```

107 def retrieve_from_hbase(self, table_name):
108     """
109     Retrieves data from HBase and returns it as a DataFrame.
110
111     Parameters:
112     - table_name: str, the name of the table to retrieve data from.
113
114     Returns:
115     - DataFrame, the retrieved data.
116     """
117     table = self.connection.table(table_name)
118     rows = table.scan()
119
120     data = []
121     for key, row in rows:
122
123         data.append(Row(
124
125             Review=row[b'cf1:Review'].decode('utf-8') if b'cf1:Review' in row else
None,
126             Sentiment=int(row[b'cf2:Sentiment'].decode('utf-8')) if
b'cf2:Sentiment' in row else None,
127             SkuInfo_index=float(row[b'cf3:SkuInfo_index'].decode('utf-8')) if
b'cf3:SkuInfo_index' in row else None,
128             tokens=row[b'cf4:tokens'].decode('utf-8').split(',') if b'cf4:tokens'
in row else None,
129             number_of_tokens=int(row[b'cf5:number_of_tokens'].decode('utf-8')) if
b'cf5:number_of_tokens' in row else None,
130         ))
131
132     schema = StructType([
133         StructField("Review", StringType(), True),
134         StructField("Sentiment", IntegerType(), True),
135         StructField("SkuInfo_index", DoubleType(), True),
136         StructField("tokens", ArrayType(StringType()), True),
137         StructField("number_of_tokens", IntegerType(), True)
138     ])
139
140
141     spark = SparkSession.builder.appName("HBase Retrieval").getOrCreate()
142     df = spark.createDataFrame(data, schema)
143
144     return df
145
146 def close(self):
147     """
148     Closes the HBase connection.
149     """
150     if self.connection is not None:
151         self.connection.close()

```