**SCHOOL OF COMPUTING (SOC)**

# CA2
# Step-by-step Tutorial

**SPECIALIST DIPLOMA IN DATA SCIENCE
(BIG DATA & STREAMING ANALYTICS)**

**IT8703 Streaming Analytics**

**Date of Submission:**     15th February 2020

# Table of Contents

# Section A – Propose your own header

## Step 1 – generate_transactions.py

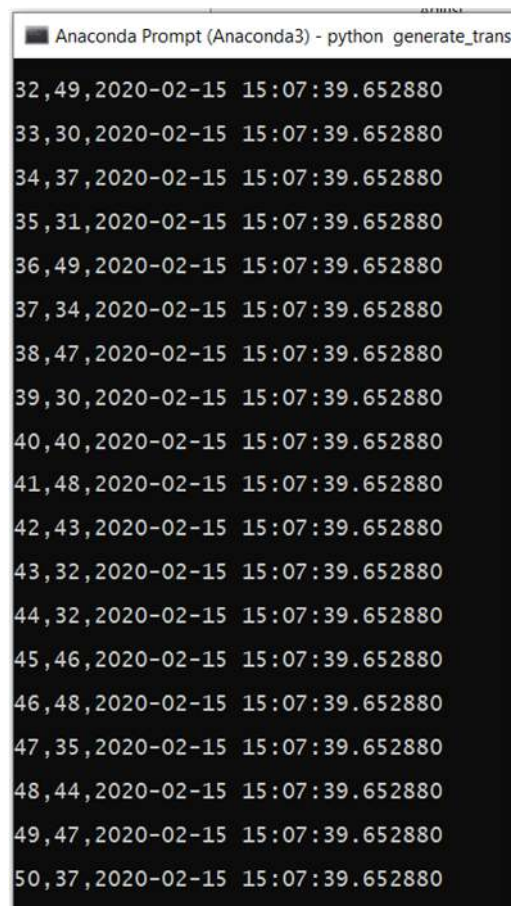| | |
|---|---|
| 1. | Random Transactions generations python code :<br><br>```python<br>from datetime import datetime<br>import sys<br>import random<br>from time import sleep<br><br><br><br># We will imagine that we have 8 rooms installed with light sensors<br># and we constantly capture the light levels of these 8 rooms every 5 seconds<br># in the subfolder iotdata<br><br># Read data<br>import pandas as pd<br>df = pd.read_csv("bakeinc_products.csv")<br><br>product_id = df['ProductID']<br><br>while True:<br>    try:<br>        dt = datetime.now()<br>        # Start a new file based on current timestamp<br>        fn = "{}{}{}{}{}{}".format(dt.year,dt.month,dt.day,dt.hour,dt.minute,dt.second)<br>        f = open("data/" + fn + ".txt","w+") # Open the file for writing<br>        for product in product_id:<br>            quantity = random.randint(30,50) # Generate a random quantity purchased<br>            data = "{},{},{}\n".format(product,quantity,dt)<br>            f.write(data) # Write into the file<br>            print(data)<br>        f.close()<br>        sleep(5) # Wait another 5 seconds before collecting the next batch of data<br>    except KeyboardInterrupt:<br>        print('Interrupted')<br>        sys.exit()<br>    except:<br>        print(sys.exc_info()[0])<br>        print(sys.exc_info()[1])<br>``` |

# Step 2 – Output of the generate_transactions.py

| | |
|---|---|
| 1. | Console output display of the transactions : |

Output the data to text file

| Name | Date modified | Type | Size |
|------|---------------|------|------|
| 202021515858.txt | 15/2/2020 3:08 pm | Text Document | 2 KB |
| 202021515853.txt | 15/2/2020 3:08 pm | Text Document | 2 KB |
| 202021515847.txt | 15/2/2020 3:08 pm | Text Document | 2 KB |
| 202021515842.txt | 15/2/2020 3:08 pm | Text Document | 2 KB |
| 202021515837.txt | 15/2/2020 3:08 pm | Text Document | 2 KB |
| 202021515831.txt | 15/2/2020 3:08 pm | Text Document | 2 KB |
| 202021515826.txt | 15/2/2020 3:08 pm | Text Document | 2 KB |
| 202021515821.txt | 15/2/2020 3:08 pm | Text Document | 2 KB |
| 202021515816.txt | 15/2/2020 3:08 pm | Text Document | 2 KB |
| 202021515810.txt | 15/2/2020 3:08 pm | Text Document | 2 KB |
| 20202151585.txt | 15/2/2020 3:08 pm | Text Document | 2 KB |
| 20202151580.txt | 15/2/2020 3:08 pm | Text Document | 2 KB |
| 202021515755.txt | 15/2/2020 3:07 pm | Text Document | 2 KB |
| 202021515749.txt | 15/2/2020 3:07 pm | Text Document | 2 KB |
| 202021515744.txt | 15/2/2020 3:07 pm | Text Document | 2 KB |
| 202021515739.txt | 15/2/2020 3:07 pm | Text Document | 2 KB |
| 2020215122613.txt | 15/2/2020 12:26 pm | Text Document | 2 KB |
| 202021512268.txt | 15/2/2020 12:26 pm | Text Document | 2 KB |
| 202021512263.txt | 15/2/2020 12:26 pm | Text Document | 2 KB |
| 2020215122558.txt | 15/2/2020 12:25 pm | Text Document | 2 KB |

# Section B – Mongo Database

## Step 1 – Run Mongod and Mongo

| 1. | Running Mongo Database Applications : |
|----|---------------------------------------|
|    | <br><br> |

# Step 2 – Saving transaction data to Mongo DB

| | |
|---|---|
| 1. | Reading streaming : CA2 - Streaming-MongoDB.ipynb<br><br>```python<br>import pandas as pd<br>import findspark<br><br>import pymongo<br>from pymongo import MongoClient<br><br>findspark.init()<br><br>from pyspark.sql import SparkSession<br>spark = SparkSession.builder.master("local[2]").appName("CA2").getOrCreate()<br><br>from pyspark.sql.types import *<br><br>filePath = "data"<br>schema = StructType().add("product_id", "integer").add("quantity", "integer").add("time", "timestamp")<br>sdf = spark.readStream.schema(schema).csv(filePath)<br><br>query = sdf.writeStream.outputMode("append").format("console")<br>query.start()<br><br><pyspark.sql.streaming.StreamingQuery at 0x276b59ee7c8><br>``` |
| 2. | Output console of the streaming :<br><br>```<br>20/02/15 15:13:18 WARN FileStreamSource: Listed 6232 file(s) in 2401.1361 ms<br>[Stage 0:>                                                         (0 + 2) / 623<br>[Stage 0:>                                                         (4 + 2) / 623<br>[Stage 0:>                                                         (20 + 2) / 623<br>[Stage 0:>                                                         (43 + 2) / 623<br>[Stage 0:>                                                         (72 + 2) / 623<br>[Stage 0:>                                                         (100 + 2) / 623<br>[Stage 0:=>                                                        (134 + 2) / 623<br>[Stage 0:=>                                                        (164 + 2) / 623<br>[Stage 0:=>                                                        (195 + 3) / 623<br>[Stage 0:=>                                                        (223 + 2) / 623<br>[Stage 0:==>                                                       (255 + 2) / 623<br>[Stage 0:==>                                                       (282 + 2) / 623<br>[Stage 0:==>                                                       (317 + 2) / 623<br>[Stage 0:===>                                                      (355 + 2) / 623<br>[Stage 0:===>                                                      (395 + 2) / 623<br>[Stage 0:===>                                                      (437 + 2) / 623<br>[Stage 0:====>                                                     (479 + 2) / 623<br>[Stage 0:====>                                                     (519 + 2) / 623<br>[Stage 0:====>                                                     (560 + 2) / 623<br>[Stage 0:=====>                                                    (600 + 2) / 623<br>[Stage 0:=====>                                                    (642 + 2) / 623<br>[Stage 0:=====>                                                    (684 + 2) / 623<br>[Stage 0:======>                                                   (726 + 2) / 623<br>[Stage 0:======>                                                   (771 + 2) / 623<br>[Stage 0:======>                                                   (815 + 2) / 623<br>[Stage 0:======>                                                   (857 + 4) / 623<br>[Stage 0:======>                                                   (901 + 2) / 623<br>[Stage 0:=======>                                                  (947 + 2) / 623<br>[Stage 0:=======>                                                  (999 + 4) / 623<br>[Stage 0:========>                                                 (1074 + 2) / 623<br>2]<br>``` |

```
-------------------------------------------
Batch: 6
-------------------------------------------
+----------+--------+--------------------+
|product_id|quantity|                time|
+----------+--------+--------------------+
|         1|      36|2020-02-15 15:15:...|
|         2|      47|2020-02-15 15:15:...|
|         3|      37|2020-02-15 15:15:...|
|         4|      33|2020-02-15 15:15:...|
|         5|      46|2020-02-15 15:15:...|
|         6|      30|2020-02-15 15:15:...|
|         7|      30|2020-02-15 15:15:...|
|         8|      41|2020-02-15 15:15:...|
|         9|      41|2020-02-15 15:15:...|
|        10|      38|2020-02-15 15:15:...|
|        11|      30|2020-02-15 15:15:...|
|        12|      50|2020-02-15 15:15:...|
|        13|      44|2020-02-15 15:15:...|
|        14|      33|2020-02-15 15:15:...|
|        15|      36|2020-02-15 15:15:...|
|        16|      41|2020-02-15 15:15:...|
|        17|      49|2020-02-15 15:15:...|
|        18|      43|2020-02-15 15:15:...|
|        19|      31|2020-02-15 15:15:...|
|        20|      37|2020-02-15 15:15:...|
+----------+--------+--------------------+
only showing top 20 rows

[Stage 12:>
```

| 3. | Saving the data to Mongo Database |
| --- | --- |

```
from pyspark.sql.functions import window,desc

window = sdf.withWatermark("time", "3 minutes").groupBy(window(sdf.time, "60 seconds")).sum('quantity')
window = window.sort(desc("window"))


def store_aggregated_data(row):

    try:
        client = MongoClient()
        # Get the database
        db = client.bakeinc
        collection = db.user

        data = {}

        time_window      = str(row["window"]["start"])
        num_transactions = int(row["sum(quantity)"])

        data['num_transactions'] = num_transactions
        data['time_window']      = time_window

        query = {'time_window': time_window}

        try:
            x = collection.find(query).next()
            collection.update_one(query,{"$set" : data})
        except StopIteration:
            collection.insert_one(data)

        print("Data Inserted")

    except KeyboardInterrupt:
        print("Keyboard Interrupted ...")
        sys.exit()

    except:
        import sys
        print("Error in store_aggregated_data")
        print(sys.exc_info()[0])
        print(sys.exc_info()[1])


agg_sdf = window.writeStream.outputMode("complete").foreach(store_aggregated_data)
agg_sdf.start()
```

```
<pyspark.sql.streaming.StreamingQuery at 0x1c4fad94688>
```

```
[Stage 12:>                (13 + 2) / 200][Stage 14:>              (0 + 0) /
2]Data Inserted
Data Inserted
[Stage 12:=>               (14 + 2) / 200][Stage 14:>              (0 + 0) /
2]Data Inserted
Data Inserted
Data Inserted
[Stage 12:=>               (15 + 2) / 200][Stage 14:>              (0 + 0) /
2]Data Inserted
Data Inserted
Data Inserted
[Stage 12:=>               (16 + 2) / 200][Stage 14:>              (0 + 0) /
2]Data Inserted
Data Inserted
[Stage 12:=>               (17 + 2) / 200][Stage 14:>              (0 + 0) /
2]Data Inserted
Data Inserted
Data Inserted
Data Inserted
Data Inserted
[Stage 12:=>               (18 + 2) / 200][Stage 14:>              (0 + 0) /
2]Data Inserted
Data Inserted
[Stage 12:=>               (19 + 2) / 200][Stage 14:>              (0 + 0) /
[Stage 12:=>               (20 + 2) / 200][Stage 14:>              (0 + 0) /
2]Data Inserted
Data Inserted
Data Inserted
Data Inserted
Data Inserted
[Stage 12:=>               (21 + 2) / 200][Stage 14:>              (0 + 0) /
2]Data Inserted
Data Inserted
[Stage 12:=>               (22 + 2) / 200][Stage 14:>              (0 + 0) /
2]
```

| | | |
|---|---|---|
| 4. | | Verification that data are inserted to Mongo Database : MongoDB_check.py |

## Reading from the bakeinc Database for verification

```
]: import pymongo
   from pymongo import MongoClient
   import pandas as pd

   client = MongoClient()
   db = client.bakeinc
   mycollections = db.user
```

```
]: data = pd.DataFrame(list(mycollections.find()))
   x = data['num_transactions']
   y = data['time_window']
```

```
]: from prettytable import PrettyTable
   table = PrettyTable(['num_transactions', 'time_window'])
   for document in db.get_collection('user').find():
       table.add_row([document['num_transactions'], document['time_window']])
   print(table)
```

```
                2844            2020-02-15 12:17:00
                2997          | 2020-02-15 12:16:00 |
                22129         | 2020-02-15 15:11:00 |
                24170         | 2020-02-15 15:10:00 |
                21807         | 2020-02-15 15:09:00 |
                4050          | 2020-02-15 15:14:00 |
                24118         | 2020-02-15 15:13:00 |
                24055         | 2020-02-15 15:12:00 |
                23762         | 2020-02-15 15:08:00 |
                7974          | 2020-02-15 15:07:00 |
                809           | 2020-02-15 12:26:00 |
                3030          | 2020-02-15 12:25:00 |
                3074          | 2020-02-15 12:24:00 |
                2993          | 2020-02-15 12:23:00 |
                2895          | 2020-02-15 12:22:00 |
                3091          | 2020-02-15 12:21:00 |
                2849          | 2020-02-15 12:20:00 |
                3103          | 2020-02-15 12:19:00 |
                2867          | 2020-02-15 12:18:00 |
   +------------------+---------------------+
```

5.

# Section C – Web-app of real-time graph

## Step 1 – Using Plotly Dash to display realtime Bar Chart

1 | Read from MongoDB and display realtime on Dash Bar Chart : MongoDB_LiveBarChart.py

```python
import dash
import dash_core_components as dcc
import dash_html_components as html
from dash.dependencies import Input, Output, State
from dash.exceptions import PreventUpdate
import random

import pymongo
from pymongo import MongoClient
import pandas as pd

external_stylesheets = ['https://codepen.io/chriddyp/pen/bWLwgP.css']

app = dash.Dash(__name__, external_stylesheets=external_stylesheets)

app.layout = html.Div(children=[
    html.H1(children='CA2 Streaming Assignment', id='first'),
    dcc.Interval(id='timer', interval=30*1000),
    html.Div(id='dummy'),
    dcc.Graph(
        id='example-graph',
    )
])

@app.callback(output=Output('example-graph', 'figure'),
        inputs=[Input('timer', 'n_intervals')])

def update_graph(n_clicks):

    client = MongoClient()
    db = client.bakeinc
    mycollections = db.user

    X = []
    Y = []

    df = pd.DataFrame(list(mycollections.find()))
```

```
          x = df['num_transactions']
          y = df['time_window']

          X = x[-20:].values
          Y = y[-20:].values

          return {
                  'data': [
                      {'x': Y,
                       'y': X,
                       'type': 'bar', 'name': 'SF'},

                  ],
                  'layout': {
                      'title': 'Dash Data Visualization'
                  }
              }


if __name__ == '__main__':
    app.run_server(debug=True)
```
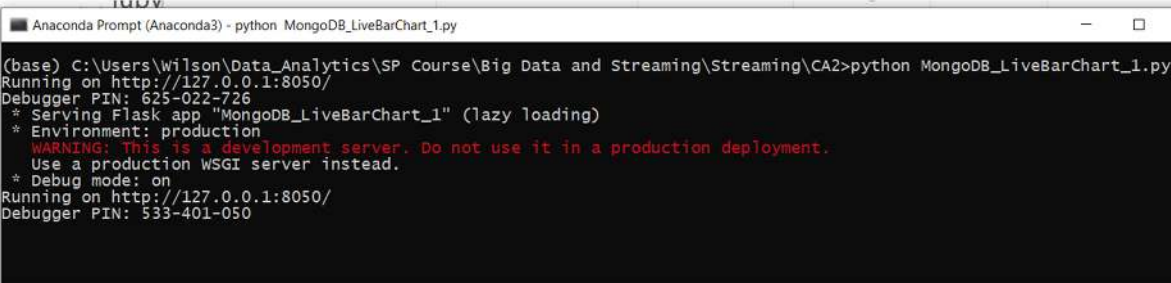
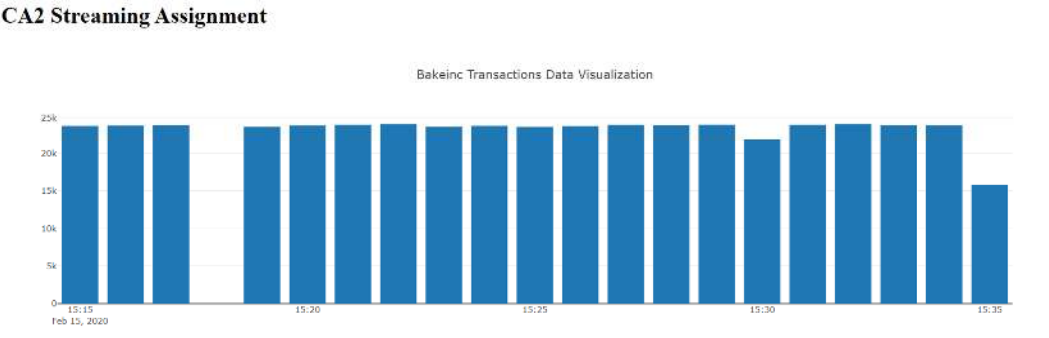| 2 | Running the dash realtime bar chart display : MongoDB_LiveBarChart_1.py |
|---|---|



| 3 | |
|---|---|

| 4 | Video submission link:<br><br>https://youtu.be/mNPKHdfJAaU |
|---|---|
| 5 | |

# Appendix A
# References

References to online materials used

**-- End of CA2 Step-by-step tutorial --**