

RNERALLA_assignment03

April 7, 2021

Assignment 3

Import libraries and define common helper functions

```
[1]: import os
import sys
import gzip
import json
from pathlib import Path
import csv

import pandas as pd
import s3fs
import pyarrow as pa
from pyarrow.json import read_json
import pyarrow.parquet as pq
import fastavro
import pygeohash
import snappy
import jsonschema
from jsonschema.exceptions import ValidationError
```

```
[8]: endpoint_url='https://storage.budsc.midwest-datascience.com'

current_dir = Path(os.getcwd()).absolute()
schema_dir = current_dir.joinpath('schemas')
results_dir = current_dir.joinpath('results')
results_dir.mkdir(parents=True, exist_ok=True)
```

```
[9]: def read_jsonl_data():
    s3 = s3fs.S3FileSystem(
        anon=True,
        client_kwargs={
            'endpoint_url': endpoint_url
        }
    )
    src_data_path = 'data/processed/openflights/routes.jsonl.gz'
    with s3.open(src_data_path, 'rb') as f_gz:
```

```

    with gzip.open(f_gz, 'rb') as f:
        records = [json.loads(line) for line in f.readlines()]

    return records

```

Load the records from <https://storage.budsc.midwest-datascience.com/data/processed/openflights/routes.jsonl.gz>

```
[10]: records = read_jsonl_data()
```

3.1 3.1.a JSON Schema

```
[11]: def validate_jsonl_data(records):
    schema_path = schema_dir.joinpath('routes-schema.json')
    validation_csv_path = results_dir.joinpath('validation-results.csv')
    with open(schema_path) as f:
        schema = json.load(f)

    with open(validation_csv_path, 'w') as f:
        for i, record in enumerate(records):
            try:
                ## TODO: Validate record
                jsonschema.validate(instance=record, schema=schema)
                pass
            except ValidationError as e:
                ## Print message if invalid record
                print("Record does not fit to the JSON schema")
                pass

```

```
[13]: validate_jsonl_data(records)
```

3.1.b Avro

```
[14]: def create_avro_dataset(records):
    schema_path = schema_dir.joinpath('routes.avsc')
    data_path = results_dir.joinpath('routes.avro')
    ## TODO: Use fastavro to create Avro dataset

    with open(schema_path) as fo:
        schema = json.loads(fo.read())
        parsed_schema = fastavro.parse_schema(schema)

    with open(data_path, 'wb') as out:
        fastavro.writer(out, parsed_schema, records)

```

```
[15]: create_avro_dataset(records)
```

3.1.c Parquet

```
[17]: def create_parquet_dataset():
    src_data_path = 'data/processed/openflights/routes.jsonl.gz'
    parquet_output_path = results_dir.joinpath('routes.parquet')
    s3 = s3fs.S3FileSystem(
        anon=True,
        client_kwargs={
            'endpoint_url': endpoint_url
        }
    )
    with s3.open(src_data_path, 'rb') as f_gz:
        with gzip.open(f_gz, 'rb') as f:
            pass
            records = [json.loads(line) for line in f.readlines()]
    df = pd.DataFrame(records)
    #with gzip.open(src_data_path, 'rb') as f:

    #print(records)
    #print(type(records))

    #with s3.open(src_data_path, 'rb') as f_gz:
    #    with gzip.open(f_gz, 'rb') as f:
    #        #pass
    #        ## TODO: Use Apache Arrow to create Parquet table and save the dataset
    table = pa.Table.from_pandas(df)
    # print(table)
    pq.write_table(table, parquet_output_path, compression='none')

create_parquet_dataset()
```

1 3.1.d Protocol Buffers

```
[19]: sys.path.insert(0, os.path.abspath('routes_pb2'))

import routes_pb2

def _airport_to_proto_obj(airport):
    obj = routes_pb2.Airport()
    if airport is None:
        return None
    if airport.get('airport_id') is None:
        return None

    obj.airport_id = airport.get('airport_id')
    if airport.get('name'):
        obj.name = airport.get('name')
```

```

if airport.get('city'):
    obj.city = airport.get('city')
if airport.get('iata'):
    obj.iata = airport.get('iata')
if airport.get('icao'):
    obj.icao = airport.get('icao')
if airport.get('altitude'):
    obj.altitude = airport.get('altitude')
if airport.get('timezone'):
    obj.timezone = airport.get('timezone')
if airport.get('dst'):
    obj.dst = airport.get('dst')
if airport.get('tz_id'):
    obj.tz_id = airport.get('tz_id')
if airport.get('type'):
    obj.type = airport.get('type')
if airport.get('source'):
    obj.source = airport.get('source')

obj.latitude = airport.get('latitude')
obj.longitude = airport.get('longitude')

return obj

def _airline_to_proto_obj(airline):
    obj = routes_pb2.Airline()
    ## TODO: Create an Airline obj using Protocol Buffers API
    if airline is None:
        return None
    if airline.get('airline_id') is None:
        return None

    obj.airline_id = airline.get('airline_id')
    if airline.get('name'):
        obj.name = airline.get('name')
    if airline.get('alias'):
        obj.alias = airline.get('alias')
    if airline.get('iata'):
        obj.iata = airline.get('iata')
    if airline.get('icao'):
        obj.icao = airline.get('icao')
    if airline.get('callsign'):
        obj.callsign = airline.get('callsign')
    if airline.get('country'):
        obj.country = airline.get('country')
    if airline.get('active'):

```

```

        obj.active = airline.get('active')
    return obj

def create_protobuf_dataset(records):
    routes = routes_pb2.Routes()
    for record in records:
        route = routes_pb2.Route()
        ## TODO: Implement the code to create the Protocol Buffers Dataset
        for key, value in record.items():
            if key=='airline':
                airline = _airline_to_proto_obj(value)
                airin = route.airline
                airin.name = airline.name
                airin.airline_id = airline.airline_id
                airin.active = airline.active
            if key=='src_airport' and value is not None:
                src_airport = _airport_to_proto_obj(value)
                srcairin = route.src_airport
                srcairin.name = src_airport.name
                srcairin.airport_id = src_airport.airport_id
                srcairin.latitude = src_airport.latitude
                srcairin.longitude = src_airport.longitude

            if key=='dst_airport' and value is not None:
                dst_airport = _airport_to_proto_obj(value)
                dstairin = route.dst_airport
                dstairin.name = dst_airport.name
                dstairin.airport_id = dst_airport.airport_id
                dstairin.latitude = dst_airport.latitude
                dstairin.longitude = dst_airport.longitude

            if key=='codeshare':
                route.codeshare = value
        routes.route.append(route)

    data_path = results_dir.joinpath('routes.pb')

    with open(data_path, 'wb') as f:
        f.write(routes.SerializeToString())

    compressed_path = results_dir.joinpath('routes.pb.snappy')

    with open(compressed_path, 'wb') as f:
        f.write(snappy.compress(routes.SerializeToString()))

create_protobuf_dataset(records)

```

3.2

3.2.a Simple Geohash Index

```
[21]: def create_hash_dirs(records):
    geoindex_dir = results_dir.joinpath('geoindex')
    geoindex_dir.mkdir(exist_ok=True, parents=True)
    hashes = []
    ## TODO: Create hash index
    for record in records:
        src_airport = record.get('src_airport', {})
        if src_airport:
            latitude = src_airport.get('latitude')
            longitude = src_airport.get('longitude')
            if latitude and longitude:
                hashes.append(pygeohash.encode(latitude, longitude))

    hashes.sort()

    three_letter = sorted(list(set([entry[:3] for entry in hashes])))

    hash_index = {value: [] for value in three_letter}

    for record in records:
        geohash = record.get('geohash')
        if geohash:
            hash_index[geohash[:3]].append(record)

    for key, values in hash_index.items():
        output_dir = geoindex_dir.joinpath(str(key[:1])).joinpath(str(key[:2]))
        output_dir.mkdir(exist_ok=True, parents=True)
        output_path = output_dir.joinpath('{}.jsonl.gz'.format(key))
        with gzip.open(output_path, 'w') as f:
            json_output = '\n'.join([json.dumps(value) for value in values])
            f.write(json_output.encode('utf-8'))
```

```
[22]: create_hash_dirs(records)
```

3.2.b Simple Search Feature

```
[23]: def airport_search(latitude, longitude):
    ## TODO: Create simple search to return nearest airport
    h = pygeohash.encode(latitude, longitude)
    dist = 0
    name = ''
    for i, record in enumerate(records):
        src_airport = record.get('src_airport', {})
        if src_airport:
```

```
lat = src_airport.get('latitude')
long = src_airport.get('longitude')
a_name = src_airport.get('name')
if lat and long:
    h1 = pygeohash.encode(lat,long)

    dist_n = pygeohash.geohash_approximate_distance(h,h1)
    if i==0:
        dist = dist_n
    else:
        if dist > dist_n:
            dist = dist_n
            name = a_name

print(name)
pass
```

```
[24]: airport_search(41.1499988, -95.91779)
```

Eppley Airfield