Task d: Streaming Application Reads stream in batches of 10 second intervals. Stores it to a separate collection in mongoDB, called "climate streaming" which is a different collection, "climate" from that in part A. Clarification • As explained the Ed forum posts, we will handle misaligned dates of hotspot and climate data by updating the hotspot data to match that of climate data which are in the SAME BATCH. • This is perfectly legal, since streaming data by its original nature, has no date and datetime information yet. • note that climate id and hotspot id are not meant to be UNIQUE within the mongo collection • however, the " id" field OR the combination of "producer" field and "climate id" or "hotspot id" field makes a particular record unique within the collection strings = ['a', 'b', 'c'] def add_comma(lst): string = "" for i in range(0, len(lst) - 1): string += f'{lst[i]}, string += str(lst[len(lst) - 1]) return string print(add comma(strings)) imports from pymongo import MongoClient
from pyspark.sql import SparkSession
from pyspark.sql.functions import col, split, element_at, when, from_json import datetime
from datetime import date from datetime import datetime from dateutil.parser import parse mongoDB stuff client = MongoClient() # connect on the default host and port db = client fit3182_assignment_db climate_streaming = db.climate_streaming # if not exist yet, then create Fieldname constants # FIELDNAMES PRODUCER = "producer" DATE = "date" LAT = "latitude" LNG = "longitude" CLIMATE_ID = "climate_id" AIR_TEMP = "air_temperature_celcius" RELATIVE_HUMIDITY = "relative_humidity" WINDSPEED = "windspeed knots" MAX_WINDSPEED = "max_wind_speed" PRECIPITATION = "precipitation" HOTSPOT_ENTRIES = "hotspot_entries" DATETIME = "datetime" HOTSPOT ID = "hotspot_id" SURFACE TEMP = "surface temperature celcius" CONFIDENCE = "confidence Other Constants HOST NAME = 'localhost' PORT NUMBER = 9092 TOPIC = "assignment" Initialize our spark session with #threads = #logicalCPU and the given application name. spark = (SparkSession builder .master('local[*]') .appName('Assignment Streaming') # localhost: 4040 getOrCreate() Create a streaming dataframe with options providing the bootstrap server(s) and topic name. kafka sdf = (spark.readStream.format('kafka') # specify the data source .option('kafka.bootstrap.servers', f'{HOST_NAME}:{PORT_NUMBER}') # can have multi .option('subscribe', f'{TOPIC}') # subscribes to the topics, can subscribe to m load() Struct Type Object StructField https://spark.apache.org/docs/latest/api/python/reference/api/pyspark.sql.types.StructField.html? highlight=structfield#pyspark.sgl.types.StructField Types https://spark.apache.org/docs/latest/api/python/reference/api/pyspark.sql.functions.from json.html? <u>highlight=from</u> <u>json</u> StructType https://spark.apache.org/docs/latest/api/python/reference/api/pyspark.sql.types.StructType.html? <u>highlight=structtype</u> $data_schema = ($ StructType() .add("producer", StringType()) .add("date", StringType()) add("climate_id", IntegerType()) add("station", StringType()) add("air_temperature_celcius", IntegerType()) add("relative_humidity", FloatType()) add("windspeed_knots", FloatType()) add("max_wind_speed", FloatType()) add("precipitation", StringType()) add("GHI_w/m2", IntegerType()) add("hotspot_id", IntegerType()) add("latitude", FloatType()) add("longitude", FloatType()) add("datetime", StringType()) add("confidence", IntegerType()) add("surface_temperature_celcius", IntegerType()) converted_stream_df = (kafka_sdf select(kafka_sdf.value.cast("string").alias("json_string") select(from_json("json_string", data_schema).alias("columns") .select("columns.*") kafka sdf In []: kafka_sdf.printSchema() converted_stream_df converted_stream_df printSchema() **Utility functions** HISTORIC PRODUCER = "climate historic" c(historic_producer = HISTORIC_PRODUCER): result = climate.delete_many({"producer": {"\$ne": historic_producer}}) return result deleted_count def format_dict_list(lst): ret = "[\n" for _dict in lst: ret += str(_dict) ret += "\n" ret += "]" return ret def format_pgh_group(groups): ret = "{\n" for key in groups: d_lst = groups[key] ret += f'"{key}": {format_dict_list(d_lst)} \n' ret += "}" return ret def format_c_dict(dic): ret = "{\n" for key in dic: if key != HOTSPOT_ENTRIES: ret += f'"{key}": {dic[key]}, ' ret += f'"{key}": {format_dict_list(dic[key])}, ' ret += "\n}' return ret def debug_print(d_nested_lst, batch_id, print_f, stage=None, stage_dataset_names = [], pr oducer_name="None"): assert_msg = f'd_nested_lst: {d_nested_lst} \n stage_dataset_names: {stage_dataset_na mes} assert(len(stage_dataset_names) == 0 or len(d_nested_lst) == len(stage_dataset_names)) if stage is not None: print('-'*20) print(f'stage: {stage}') print('-'*20) elif stage is None: print('-'*20) print(f"batch: {batch_id}") print('-'*20) Len(stage_dataset_names) == 0: for dict_lst in d_nested_lst: producer = dict_lst[0]["producer"] if len(dict_lst) > 0 else producer_name print(f'{len(dict_lst)}, producer: {producer}, data: {print_f(dict_lst)}') elif len(stage dataset names) > 0 and len(stage dataset names) == len(d nested lst): for i in range(len(d_nested_lst)): dict_lst = d_nested_lst[i] string = f'{len(dict_lst)}, string += f'stage_dataset: {stage_dataset_names[i]}, ' string += f' data: {print_f(dict_lst)} print(string) errmsg = f'd_nested_lst: {d_nested_lst} \n stage_dataset_names: {stage_datase} t names}' raise IndexError(errmsg) errmsg = f'd_nested_lst: {d_nested_lst} \n stage_dataset_names: {stage_dataset_na mes}' len1, len2 = ler(stage_dataset_names), ler(d_nested_lst) raise Exception(f"len(stage_dataset_names) {len1} != len(d_nested_lst) {len2}, , {errmsq}") def separate_producer(lst, producer_name_lst): res = {} # list of dictionaries, key is producer name for producer in producer_name_lst: res[producer] = [] for dict_obj in lst: producer = dict_obj["producer"] res[producer] append(dict_obj) return res def str2date(x, field_name): if x[field_name] is not None: parsed = parse(x[field_name], ignoretz=llone, dayfirst=True) $x[field_name] = parsed$ except TypeError: errmsg = f'{x}, {field_name} raise TypeError(errmsg) def my_sum(lst): res = 0 for x in lst: res += x return res def my_avg(lst, intFlag = True): int_avg = my_sum(lst) / lem(lst) if intFlag: int_avg = round(int_avg) return int_avg except TypeError: errmsg = str(lst)raise TypeError(errmsg) same_pgh_code(d_lst, lat_lng_to_check, precision_val): pgh_check_code = pgh_encode(*lat_lng_to_check, precision = precision_val) for dic in d_lst: code = pgh.encode(dic[LAT], dic[LNG], precision=precision_val) if pgh_check_code != code: def check all same pgh_code(d_lst, lat_lng_to_check, precision_val): pgh_check_code = pgh.encode(*lat_lng_to_check, precision = precision_val) for dic in d_lst: code = pgh_encode(dic[LAT], dic[LNG], precision=precision_val) if pgh_check_code != code: def pgh_from_pair(pair, precision_val): return pgh.encode(*pair, precision = precision_val) ode(pair1, pair2, precision_val): code1 = pgh_from_pair(pair1, precision_val) code2 = pgh_from_pair(pair2, precision_val) return code1 == code2 def lat_lng_pair_from_dict(dic): return dic[LAT], dic[LNG] def sort_dict_lst_on_field(d_lst, field_name): f = lambda dic : dic[field_name] return sorted(d_lst, f) gh_grouping(d_lst): groups = {} for rec in d_lst: pgh_code = pgh.encode(rec[LAT], rec[LNG], precision=5) if pgh_code not in groups: groups[pgh_code] = [rec] groups[pgh_code].append(rec) return groups def agg_pgh_groups(pgh_groups, field_name_lst, agg_f_lst, producer_lst): assert(len(field_name_lst) == len(agg_f_lst)) agg_groups = copy.deepcopy(pgh_groups) # make a copy first to prevent mutations for pgh_key in agg_groups: group = agg_groups[pgh_key] if len(group) > 1: agg_rec = copy.deepcopy(group[0]) # initialize to first record's copy for i in range(len(field_name_lst)): print(f"agg_pgh_groups: ran agg loop {i+1} times !") field, agg_f = field_name_lst[i], agg_f_lst[i] field_values = [dic[field] for dic in group] # extract fields value agg_val = agg_f(field_values) agg_rec[field] = agg_val # update value in agg record agg_rec[PRODUCER] = f'merged: {producer_lst}' # update producer to indicate m print(f"agg_rec: {agg_rec} !") agg_groups[pgh_key] = [agg_rec] print(f"updated groups with agg_rec: {agg_groups} !") return agg_groups l(nested): ret = [] for lst in nested: ret += lst return ret Db_writer and Console_logger BATCH_PROCESSING_TIME = 10 def batchFunc(batch_df, batch_id): raw_data = batch_df.collect() # returns a list of Row objects dict_lst = [row.asDict() for row in raw_data] # PREPROCESSING climate_data = [x for x in dict_lst if x["climate_id"] is not None] hotspot_data = [x for x in dict_lst if x["hotspot_id"] is not None] debug_print([climate_data, hotspot_data], batch_id, stage=None, print_f = format_dict_list, stage_dataset_names = ["climate_data", "hotspot_data"]) len(climate_data) > 0: climate_data = [str2date(x, "date") for x in climate_data] if len(hotspot_data) > 0: hotspot_data = [str2date(x, "date") for x in hotspot_data] hotspot_data = [str2date(x, "datetime") for x in hotspot_data] debug_print([climate_data, hotspot_data], batch_id, print_f = format_dict_list, stage=f"after convert str to datetime", stage_dataset_names = ["climate_data", "hotspot_data"]) if len(climate data) == 0: debug_print([climate_data], batch_id, print_f = format_dict_list, stage=f"no climate data") elif len(climate_data) > 1: raise ValueError("there is more than 1 climate data in this batch!") elif len(climate data) == 1: c_dict = climate_data[0] climate_date = c_dict["date"] for row dict in hotspot data: date_time = row_dict['datetime'] row_dict["date"] = climate_date year = climate_date.year month = climate_date.month day = climate_date.day hour, minute, second = date_time.hour, date_time.minute, date_time.second new_datetime = datetime(year, month, day, hour, minute, second) row_dict['datetime'] = new_datetime hotspot_aqua_and_terra = separate_producer(hotspot_data, ["hotspot_AQUA", "hotspo t_TERRA"]) hotspot_aqua = hotspot_aqua_and_terra["hotspot_AQUA"] hotspot_terra = hotspot_aqua_and_terra["hotspot_TERRA"] debug_print([climate_data, hotspot_aqua, hotspot_terra], batch_id, print_f = format_dict_list, stage = "updated hotspot date and datetime to match climate's date") c_dict = climate_data[0] c_geohash = pgh.encode(c_dict[LAT], c_dict[LNG], precision=3) loc_filtered_hotspots = [] for h_dict in hotspot_data: h_geohash = pgh.encode(h_dict[LAT], h_dict[LNG], precision=3) if h_geohash == c_geohash: loc_filtered_hotspots append(h_dict) debug_print([loc_filtered_hotspots], batch_id, print_f = format_dict_list, stage=f"hotspot loc filter {c_dict[LAT]}, {c_dict[LNG]}", stage_dataset_names=["loc_filtered_hotspots"]) if len(loc_filtered_hotspots) == 0: $embedded = c_dict$ embedded["hotspot_entries"] = [] climate_streaming insert_one(embedded) debug_print([climate_data, loc_filtered_hotspots], batch_id, stage= f"only climate data after loc_filter, inserted to db", print_f = format_dict_list, stage_dataset_names = ["climate_data", "loc_filtered_hotspots"]) elif len(loc_filtered_hotspots) > 0: pgh_groups = pgh_grouping(loc_filtered_hotspots) debug_print([pgh_groups], batch_id, print_f = format_pgh_group, stage=f"after groups, before agg", stage_dataset_names = ["pgh_groups"]) agg_groups = agg_pgh_groups(pgh_groups, [SURFACE_TEMP, CONFIDENCE], [my_avg, my_avg], ["AQUA", "TERRA"]) debug_print([agg_groups], batch_id, print_f = format_pgh_group, stage=f"groups after aggregate", stage_dataset_names = ["aggregated_h_groups"]) agg_groups = [agg_groups[pgh_key] for pgh_key in agg_groups] agg_hotspot = unpack_all(agg_groups) natural_fire_flag = c_dict["air_temperature_celcius"] > 20 and c_dict["GHI_w/ fire_cause = "natural" if natural_fire_flag else "other" for h_dict in agg_hotspot: h_dict["fire_cause"] = fire_cause c dict = climate_data[0] c_dict["hotspot_entries"] = agg_hotspot debug_print([agg_hotspot], batch_id, print_f = format_dict_list, stage=f"final", stage_dataset_names = ["finalized data"]) climate_streaming insert_one(c_dict) debug_print([c_dict], batch_id, print_f = format_c_dict, stage=f"inserted to db", stage_dataset_names = ["c_dict"]) db_writer = (converted_stream_df writeStream .outputMode('append') .trigger(processingTime = f'{BATCH_PROCESSING_TIME} seconds') . foreachBatch(batchFunc) db_writer start streaming query writer = db_writer query = writer.start() query awaitTermination() except StreamingQueryException as exc: print(exc) query stop()