

Big Data Analytics Assignment

Name : Singam Adikesava Reddy

VTU NO : VTU25083

course code : 10211IT201

Subject : Big Data Analytics

Task No : (3)

JK 27/12/15

(5/5)

Dimensions

dim_time

- time_id (int)
- date (Date)
- day-of-week, day-of-month, week-of-year
- is-holiday Boolean
- fiscal-month, fiscal-quarter, fiscal-year

dim_store

- Store_id (int)
- Store-code, name
- region, city, state, country
- store-type (flagship, outlet)
- Store-open-date, store-close-date
- Square-footage, manager-id

dim_Product

- Product_id (int)
- SKU, UPC
- name, brand, category, sub-category
- size, color
- sale-start-date, sale-end-date
- current-price, cost-price
- product-status (active | discontinued)

dim_customer

- customer_id (int)
- customer-type (quest)
- first-name, last-name
- gender, birth-year
- loyalty-tier, join-date

```
from pyspark import sparkContext
from pyspark.streaming import StreamingContext
from pyspark.streaming.kafka import KafkaUtils
import json
```

```
ssc = sparkContext.appName = "video"
ssc = StreamingContext(ssc, batchDuration=5)
ssc.checkpoint("/path/checkpoint")
kvs = KafkaUtils.createDirectStream(ssc, topics
["video-events"], kafkaParams={"metadata.broker.list":
"192.168.1.1:9092", "zookeeper.connect": "192.168.1.1:2181"})
plays = events.map(lambda msg: json.loads(msg[1]))
["play-start", "play", "play-end"])

content = plays.map(lambda e: (e["content-id"], e.get("play-seconds", 0)))
batch = plays.map(lambda e: (e["content-id"], e.get("play-seconds", 0)))

By key (lambda a,b: a+b)
```

```
def update(running, total):
    return sum(newVals) + (running - total)

running = total = 0
ssc.start()
```

```
ssc.awaitTermination()
```

2) OLAP data Model for xyz Retail (Sales Analytics)
support fast multidimensional analysis of sales across,
store, time, products, customers, channel and
Promotions (kpls: sales amount)

Recommended model

A star schema with one central fact table and
several conformed dimension tables

Fact - Sales

surrogate key: Sales-id (int)

fk-time-id (int)

fk-store-id (int)

fk-store-product-id (int)

fk-customer-id (int)

fk-channel-id (int)

fk-promo-id (int)

quantity-sold (int)

Sales-amount (DECIMAL(12,2))

unit-price (DECIMAL(10,2))

discount-amount (DECIMAL(12,2))

returned-flag (Boolean)

transaction-type (string)

created-at (time stamp)

load-batch-id (varchar)

partition this fact time (Year | month)

Implementing Apache Ignite for AlphaBank Goals

3.
 - Speed up transaction and analytical operations via in-memory data grid
 - Provide low latency caching for customer profiles account balances reference data and accelerate OLTP
 - Use Ignite for distributed in-memory SQL, key-value, compute grid learning

High-level architecture

- Data tier: persistent store (RDBMS/postgreSQL) + Ignite cluster as in-memory layer with optional Ignite cluster: multiple node across data centers (replication, partitioning)
- client: banking application server
- integration: banking application use Ignite cache store for read-through/write-through RDBMS.

Implementation steps

1. select deployment mode
 - option A: Ignite as distributed cache in front of RDBMS (recommended for minimal change)
 - Option B: Ignite with native persistence to replace some DB workloads.

2. Define data model

accounts, customer, transactions - recent,

product - rates

3. configure cache store

, implement cache store interface to sync with

PDMS, enabling persistence backing.

, configure write behind for better throughput

put (tunable).

4. set up clustering

, Data partitioned cache replicated

select c.customer_id, account_id, a.balance
from accounts a.

Join customers (on a.customer_id = c.

customer_id

where c.region = 'south';

Implement Apache Hive for Retail Groove Objectives

- 1) Centralize Customer & transaction data for large-scale analytics marketing
- 2) Provide SQL-like interface (HiveQL) over data lakes

customer table

customer_id BIGINT,

first_name string

last_name string

email_hash string

sign-up date Date,

loyalty-tier string

gender string

birth-year int,

stored AS Parquet

location 'data / retail / customers';

Transactions table

create external Table retail_transactions

transaction_id, string,

customer_id BIGINT,

store_id int,

product_sku string

quantity int

unit_price Decimal(10, 2))

1. Monthly revenue trend for a Product category
Select t-year, t-month, sum(f.sales-amount)
As revenue from fact-sales f.

Join dim-product P ON f.fk-product-id =
P.Product-id

Join dim-time t ON t.fk-time-id = t.time-id

Where P.category = 'electronics' AND t.date >=
DATEADD(month, -12, current-date)

GROUP BY t-year, t-month

ORDER BY t-year, t-month;

2. Top 10 SKUs by margin last quarter

Select p.sku, p.name, s.region, sum(f.margin-
amount) As total-margin
From fact-sales f.

Join dim-product P ON f.fk-product-id = P.Product-id

Join dim-store S ON f.fk-store-id = S.store-id

Join dim-time t ON f.fk-time-id = t.time-id.

Where t.fiscal-quarter = 3 AND t.fiscal-year
= 2025

GROUP BY P.sku, P.name, S.region

ORDER BY total-margin Desc.

Limit 10;

Product table

Create external table retail-products (

SKU string,

name string, brand string, category string,

Price Decimal (10,2), cost Decimal (10,2),

cost Decimal (10,2), active Boolean)

stored AS Parquet

location '/data/retail/products/';

Predict the future of Hadoop

Key directions

1. Move from HDFS - centric to object-store - first
2. Hadoop ecosystems will continue to integrate with cloud object stores (S3, GCS) HDFS will remain for on-prem but the trend favors object store connector.
3. Smaller foot print, more modular
4. Monolithic Hadoop distributions will give way to modular components (Yarn alternative)
5. Convergence towards SQL
6. Data formats (Parquet/ORC) table formats (Apache Iceberg, Delta Lake) stack.
7. Stream + batch unification
8. Real time processing engine (Flink, Kafka Streams, Spark structured streaming) retired.
9. Cloud-native, containerized
10. Kubernetes orchestration for data workloads, auto-scaling, improved resource operator patterns for components like Hive, HDFS alternatives.

7. integration with ML/AI workflows

- Native support for model training integration

8. performance & cost efficiency

- More compute separation from storage with caching layer (Alluxio, Ignite), tiered storage strategies (hot in memory, warm SSD, cold object store).

9. serverless analytics

- More serverless SQL-on-lake offerings (Athena, BigQuery like) reduce operator overhead, but Hadoop tools will adapt to expose serverless capabilities.