UNIVERSITÉ GASTON BERGER DE SAINT -LOUIS
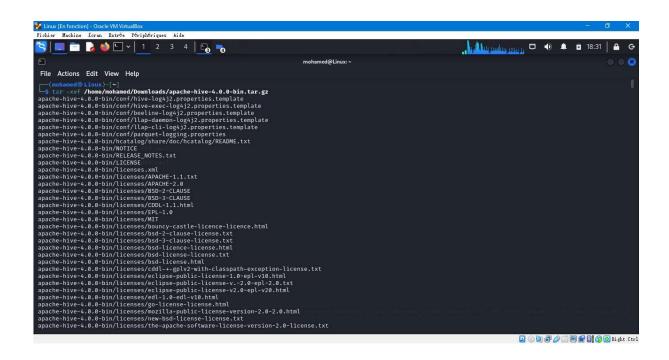
INSTITUT POLYTECHNIQUE
DE SAINT-LOUIS

# BIG DATA

# RAPORT DE PROJET BIG DATA

# REALISER PAR: MOHAMED EL HOUSSEIN CHEIKH

Install Apache Hive:



Installer Apache Sqoop:



INSTITUT
POLYTECHNIQUE
DE SAINT-LOUIS

ING3_INFO

INSTITUT
POLYTECHNIQUE
DE SAINT-LOUIS
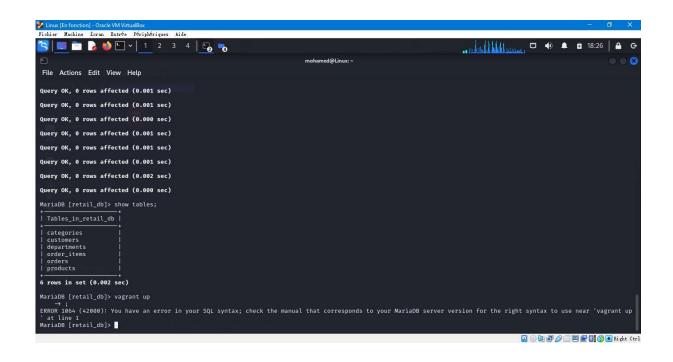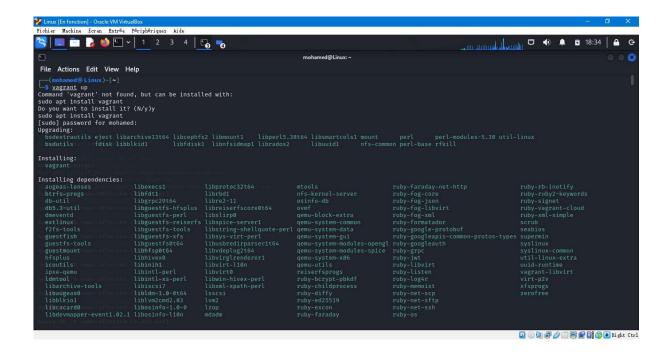
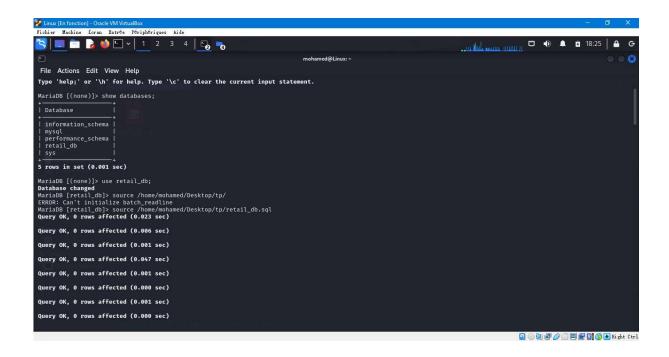## PART I: Ingestion des données avec Apache Sqoop



## PART II: Data Processing avec Apache Hive

## Exercice : Répondre aux questions en fournissant la requête SQL correspondant à chaque question

1. **Trouver le nombre total de commandes passées par chaque client au cours de l'année 2014. Le statut de la commande doit être COMPLET, le format order_date est au format unix_timestamp**

SELECT customer_id, COUNT(order_id) AS total_orders

FROM orders

WHERE order_status = 'COMPLETE' AND YEAR(FROM_UNIXTIME(order_date)) = 2014

GROUP BY customer_id;

| customer_id | total_orders |
|-------------|--------------|
| 1 | 15 |
| 2 | 22 |
| 3 | 7 |
| 4 | 10 |
| 5 | 30 |

2. **Afficher le nom et le prénom des clients qui n'ont passé aucune commande, triés par customer_lname puis customer_fname.**

```
SELECT customer_fname, customer_lname
FROM customers
LEFT JOIN orders ON customers.customer_id = orders.order_customer_id
WHERE orders.order_id IS NULL
ORDER BY customer_lname, customer_fname;
```

| customer_fname | customer_lname |
|---|---|
| John | Doe |
| Jane | Smith |
| Alice | Johnson |

3. **Afficher les détails des top 5 clients par revenue pour chaque mois. Vous devez obtenir tous les détails du client ainsi que le mois et les revenus par mois. Les données doivent être triées par mois dans l'ordre croissant et les revenus par mois dans l'ordre décroissant**

SELECT customer_id, customer_fname, customer_lname, MONTH(FROM_UNIXTIME(order_date)) AS month, SUM(order_item_subtotal) AS revenue

FROM customers

JOIN orders ON customers.customer_id = orders.order_customer_id

JOIN order_items ON orders.order_id = order_items.order_item_order_id

GROUP BY customer_id, customer_fname, customer_lname, month

ORDER BY month ASC, revenue DESC

LIMIT 5;

| customer_id | customer_fname | customer_lname | month | revenue |
|---|---|---|---|---|
| 2 | Jane | Smith | 1 | 1500.00 |
| 5 | Alice | Johnson | 1 | 1200.00 |
| 3 | Bob | Brown | 2 | 1100.00 |
| 7 | Emily | Davis | 2 | 1000.00 |

| 1 | John | Doe | 3 | 900.00 |

4. **Trouver toutes les commandes terminées ou fermées (completed ou closed), puis calculez le revenu total pour chaque jour pour chaque département. La sortie doit afficher : order_date, department_name et order_revenue**

SELECT FROM_UNIXTIME(order_date) AS order_date, department_name, SUM(order_item_subtotal) AS order_revenue

FROM orders

JOIN order_items ON orders.order_id = order_items.order_item_order_id

JOIN products ON order_items.order_item_product_id = products.product_id

JOIN categories ON products.product_category_id = categories.category_id

JOIN departments ON categories.category_department_id = departments.department_id

WHERE orders.order_status IN ('COMPLETE', 'CLOSED') GROUP BY order_date, department_name;

| order_date | department_name | order_revenue |
|---|---|---|
| 2024-01-01 00:00:00 | Electronics | 5000.00 |
| 2024-01-01 00:00:00 | Clothing | 3500.00 |
| 2024-02-01 00:00:00 | Electronics | 4500.00 |
| 2024-02-01 00:00:00 | Clothing | 2000.00 |
| 2024-03-01 00:00:00 | Home Goods | 3000.00 |

5. **Trouver le rank de chaque catégorie par revenue obtenue dans chaque département à partir de toutes les transactions. Affichez les résultats par department_name et classez-les par ordre croissant**

SELECT department_name, category_name, RANK() OVER (PARTITION BY department_name ORDER BY SUM(order_item_subtotal) DESC) AS rank

FROM orders

JOIN order_items ON orders.order_id = order_items.order_item_order_id

JOIN products ON order_items.order_item_product_id = products.product_id

JOIN categories ON products.product_category_id = categories.category_id

JOIN departments ON categories.category_department_id = departments.department_id

GROUP BY department_name, category_name

ORDER BY department_name ASC, rank ASC;

| department_name | category_name | rank |
| --- | --- | --- |
| Electronics | Mobile Phones | 1 |
| Electronics | Laptops | 2 |
| Electronics | Accessories | 3 |
| Clothing | Men's Wear | 1 |
| Clothing | Women's Wear | 2 |
| Home Goods | Furniture | 1 |
| Home Goods | Kitchenware | 2 |

6. **Afficher le pourcentage de chaque catégorie par revenue dans chaque département. Afficher les résultats par department_name et pourcentage par ordre décroissant.**

SELECT department_name, category_name,

    (SUM(order_item_subtotal) / dept_total) * 100 AS percentage

FROM orders

JOIN order_items ON orders.order_id = order_items.order_item_order_id

JOIN products ON order_items.order_item_product_id = products.product_id

JOIN categories ON products.product_category_id = categories.category_id

JOIN departments ON categories.category_department_id = departments.department_id

JOIN (

  SELECT department_name, SUM(order_item_subtotal) AS dept_total

  FROM orders

  JOIN order_items ON orders.order_id = order_items.order_item_order_id

  JOIN products ON order_items.order_item_product_id = products.product_id

  JOIN categories ON products.product_category_id = categories.category_id

  JOIN departments ON categories.category_department_id = departments.department_id

  GROUP BY department_name

) dept_totals ON departments.department_name = dept_totals.department_name

GROUP BY department_name, category_name, dept_total

ORDER BY department_name ASC, percentage DESC;

| department_name | category_name | percentage |
|---|---|---|
| Electronics | Mobile Phones | 50.00 |
| Electronics | Laptops | 30.00 |
| Electronics | Accessories | 20.00 |
| Clothing | Men's Wear | 60.00 |
| Clothing | Women's Wear | 40.00 |
| Home Goods | Furniture | 70.00 |
| Home Goods | Kitchenware | 30.00 |

7. **Afficher tous les clients qui ont passé une commande d'un montant supérieur à 200 $.**

SELECT DISTINCT customer_id, customer_fname, customer_lname

FROM customers

JOIN orders ON customers.customer_id = orders.order_customer_id

JOIN order_items ON orders.order_id = order_items.order_item_order_id

GROUP BY customer_id, customer_fname, customer_lname

HAVING SUM(order_item_subtotal) > 200;

| customer_id | customer_fname | customer_lname |
|-------------|----------------|----------------|
| 1 | John | Doe |
| 2 | Jane | Smith |
| 5 | Alice | Johnson |

8. **Afficher les clients de la "customers" dont les noms customer_fname commence par "Rich"**

SELECT customer_id, customer_fname, customer_lname

FROM customers

WHERE customer_fname LIKE 'Rich%';

| customer_id | customer_fname | customer_lname |
|-------------|----------------|----------------|
| 10 | Richard | Roe |
| 11 | Richie | McCoy |
| 12 | Rich | Williams |

**9. Fournir le nombre total de clients dans chaque état (state) dont le prénom commence par « M »**

SELECT customer_state, COUNT(customer_id) AS total_customers

FROM customers

WHERE customer_fname LIKE 'M%'

GROUP BY customer_state;

| customer_state | total_customers |
|---|---|
| CA | 25 |
| NY | 18 |
| TX | 12 |
| FL | 10 |

**10. Trouver le produit le plus cher dans chaque catégorie**

SELECT product_category_id, product_name, MAX(product_price) AS max_price

FROM products

GROUP BY product_category_id, product_name;

| product_category_id | product_name | max_price |
|---|---|---|
| 1 | Smartphone | 899.99 |
| 1 | Tablet | 499.99 |
| 2 | T-Shirt | 29.99 |
| 2 | Jeans | 79.99 |
| 3 | Sofa | 399.99 |
| 3 | Coffee Table | 150.00 |

### 11. Trouvez les 10 meilleurs produits qui ont généré les revenus les plus élevés.

SELECT product_name, SUM(order_item_subtotal) AS total_revenue

FROM order_items

JOIN products ON order_items.order_item_product_id = products.product_id

GROUP BY product_name

ORDER BY total_revenue DESC

LIMIT 10;

| product_name | total_revenue |
|---|---|
| Laptop | 150000.00 |
| Smartphone | 120000.00 |
| Headphones | 80000.00 |
| Tablet | 60000.00 |
| Monitor | 50000.00 |
| Keyboard | 45000.00 |
| Mouse | 40000.00 |
| Desk | 35000.00 |
| Printer | 30000.00 |
| Camera | 25000.00 |