

Challenges Big data - 1 BIA 2022-2023



Encadré par :

➤ Pr. El HAJJI Mohamed

Réalisé par :

➤ BENAABIDATE Ahmed Walid.

Challenges

Big data - 1

BIA 2022-2023

Challenge 1

Proposer une solution pour mettre en œuvre un système d'analyse de sentiments basé sur une solution Big Data basée sur les commentaires des visiteurs de site web HESPRESS, qui supporte le data Ingestion, le streaming, le batch processing, et le dashboarding temps réel.

Solution

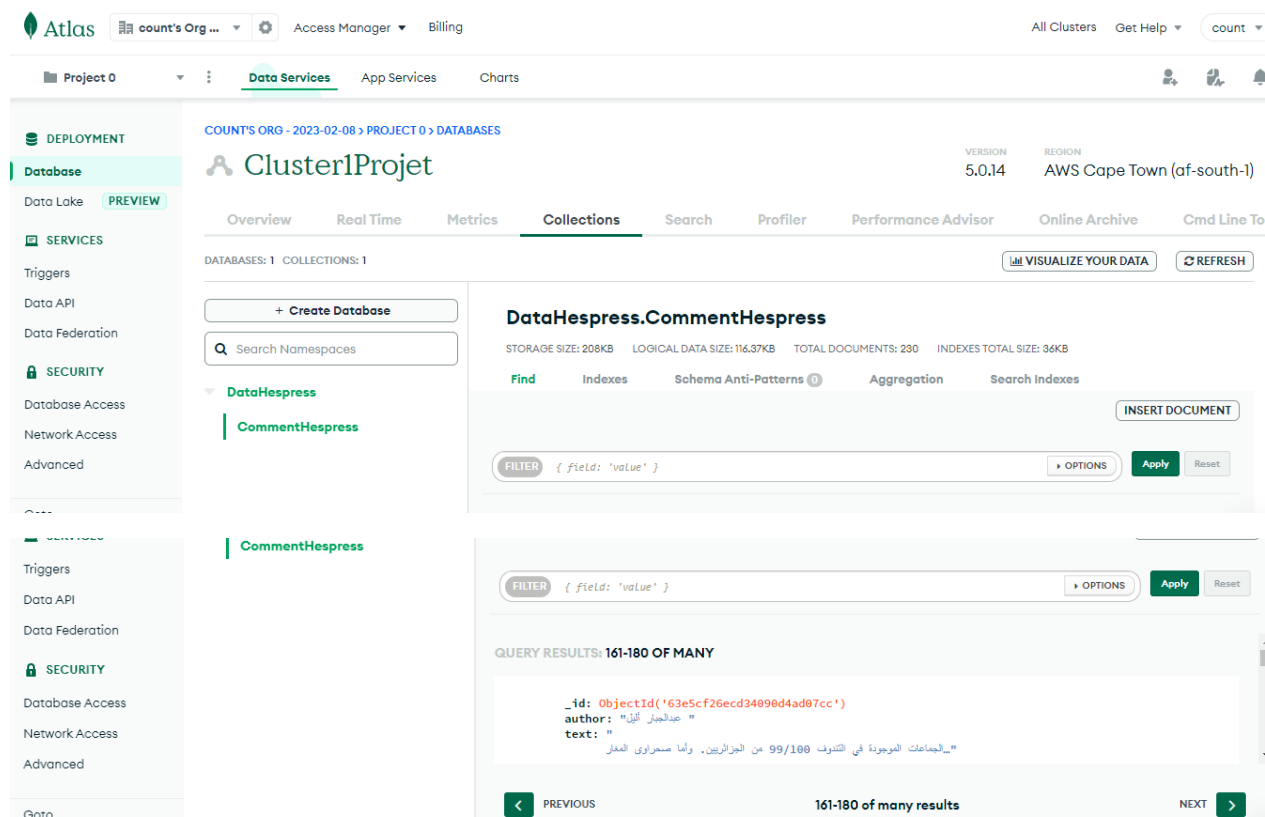
Pour mettre en œuvre un système d'analyse de sentiments pour les commentaires des visiteurs du site web HESPRESS, on a choisi une architecture qui permet de traiter et de visualiser les commentaires des visiteurs en temps réel pour déterminer leur sentiment.

Voilà l'architecture :

- + Data Ingestion: Configurer Apache Kafka pour ingérer en temps réel les commentaires des visiteurs dans un cluster distribué.
- + Streaming: Utiliser Apache Spark pour traiter les commentaires en temps réel, en utilisant les flux de données provenant d'Apache Kafka, et en les analysant pour déterminer le sentiment.
- + Batch Processing: Configurer Apache Cassandra pour stocker les données des commentaires en grande quantité, ce qui permettra de traiter les données en lots pour obtenir des insights plus approfondis sur les tendances de sentiment au fil du temps.
- + Dashboarding temps réel: Utiliser un outil de tableau de bord en temps réel pour visualiser les résultats de l'analyse de sentiments.

D'abord on a commencé par le Scraping des commentaires du site web **HEPRESS**

On a créé une base de données sur « **mongoDB** » pour stocker les commentaires



Récupération et stockage des commentaires de HESPRESS via l'utilisation de la technique de web Scraping dans MongoDB.

```
import requests
from bs4 import BeautifulSoup
from pymongo import MongoClient

# Définir L'URL à scraper
url = "https://www.hespress.com/%d8%a5%d9%86%d9%87%d8%a7%d8%a1-%d9%85%d9%87%d8%a7%d9%85-%d8%b3%d9%81%d9%8a%d8%b1-%d8%a7%d9%84%d9%85%d8%ba%d8%b1%d8%a8-%d8%a8%d9%81%d8%b1%d9%"
# Faire une requête GET au site Web
response = requests.get(url)

# Analyser HTML en utilisant BeautifulSoup
soup = BeautifulSoup(response.text, "html.parser")

# Extraire Les commentaires du HTML
comments = soup.find_all("div", class_="comments")

# print(comments)

# Connexion à l'instance MongoDB
client = MongoClient("mongodb+srv://walid:123@cluster1projet.84eixxu.mongodb.net/?retryWrites=true&w=majority")

# Obtenir La collection de commentaires
comments_collection = client.DataHespress.CommentHespress

comments_to_insert = []

# Parcourez Les commentaires et extrairez Le texte des deux éléments div
for comment in comments:
    for ul in comment.find_all('ul'):
        for li in ul.find_all('li'):
            div_container = li.find('div', {'class': 'comment-body'})
            div1 = div_container.find('div', {'class': 'comment-head'})
            div2 = div_container.find('div', {'class': 'comment-text'})
            comment_obj = {"author": div1.text.split("\n")[2], "text": div2.text}
            comments_to_insert.append(comment_obj)

# insérer tous Les commentaires
comments_collection.insert_many(comments_to_insert)
```

Les commentaires extraits de la page Web HESPRESS sont ingérés dans un topic Kafka appelé « **topicExam** »

[illegible]

On utilise Spark pour traiter les données par lots en groupant les commentaires selon l'auteur et en transformant les valeurs en une liste.

```
from pyspark.sql import SparkSession

# Créer une session Spark
spark = SparkSession.builder.appName("BatchProcessing").getOrCreate()

# Charger Les données du fichier JSON Local dans un DataFrame
df = spark.read.option("multiline", "true").json("c:\Users\hp2\Desktop\CommentHespress.json")

# Effectuer un traitement par lots sur Les données
result = df.rdd\
    .map(lambda x: (x.author, x.text))\
    .groupByKey()\
    .mapValues(list)\
    .collect()

# Show result
for author, text in result:
    print(f"Author: {author}\nText: {text}\n")

# Afficher Le résultat
spark.stop()
```

Author: بولخن
 Text: أغلبية العقليات المغرب خاسر وأوروبا رابحة. لا توجد العقلاية رابح راجح. المغرب اجر سواخله للتواجر الأوربي ب 30مليون يورو على مدى خمس سنوات. في الوقت الذي تريح أوروبا ١٠٠مليون يورو على مدى 30 سنة. انه اكبر نيب لمقدرات المغاربة. 30 مليون يورو لا تأتي حتى لمبلغ او مستثنى متوسط او طريق بطول 5كم [١٠]. 4. ملايين يزور من سلكه المغرب سنويا. انما اكبر نيب لمقدرات المغاربة.

هذه من يعطي صورة جميلة عن بلدته وهذه من يشترى العافيا المغربية في هولندا . التفتل الفساد القتل الجريمة المنظمة. المخدرات. ولهذا وهولندا تريد إرجاع هؤلاء المجرمين إلى وطنهم الأم. لاجل ولا قوة إلا بالله

Text: ['n']
 المحرقات عصب الحياة . لهذا من الصعب أن هناك شائقة وسطاء في كل شيء. واحد يبيع لآخر وهكذا. الأمم الضعيف على شركات المحرقات لإرجاع من المحرقات إلى ثمنه الأصلي 'n' . بها وبخشيها بأعمال إجرامية

['n'] . سف لو كنا شعبا واعيا لفلاننا محطة بعلينا لمدة طويلة ولجئنا يوم الأحد يوم بدون محطة أو بدون سيارة لعرض مداخل المحطات من جهة و للحد من التلوث من جهة أخرى

Author: عبدالله 768
Text: ["\nالمادّا تحكروون المغرب\nمن قبل لم تكن اتفاقية لتبادل المجرمين اما اذا تم توكلّمها فسرى هل سيتم تسليمه ام لا ما عدا اذا مهدت له هولندا الطريق للهرب ليلد آخر\nاردا على الذي يدافع عن شعور\n"]

On utilise Spark pour effectuer un traitement par lots sur les données en appliquant une analyse des sentiments à l'aide de la bibliothèque « **TextBlob** »

```
from pyspark.sql import SparkSession
from textblob import TextBlob

# Créer une session Spark
spark = SparkSession.builder.appName("BatchProcessing").getOrCreate()

# Charger Les données du fichier JSON Local dans un DataFrame
df = spark.read.option("multiline", "true").json(r"C:\Users\hp2\Desktop\CommentHespress.json")

# Effectuer un traitement par lots sur Les données
result = df.rdd\
    .map(lambda x: (x.author, x.text))\
    .groupByKey()\
    .mapValues(list)\
    .collect()

# Effectuez une analyse des sentiments à l'aide de TextBlob
for author, text in result:
    sentiment_scores = [TextBlob(t).sentiment.polarity for t in text]
    average_sentiment = sum(sentiment_scores) / len(sentiment_scores)

    print(f"Author: {author}\nAverage sentiment: {average_sentiment}\n")

# Arrêter La session Spark
spark.stop()
```

Author: أحمد
Average sentiment: 0.0

Author: العريشي
Average sentiment: -0.75

Author: زكرياء
Average sentiment: 0.0

Author: صفريوي
Average sentiment: 0.0

Author: kala
Average sentiment: 0.0