Projet Base de données NoSQL



Présenté par : Astrid Aurelien NKUMBE ENONGENE Anne-Josée LOUIS

Sommaire:

Sommaire :	
1. Introduction	2
Contexte	2
Objectifs	2
Technologies Utilisées	2
2. Architecture Système	3
Détail des Composants	3
Interaction entre les Composants	4
3. Modèle de Données	4
MongoDB	4
Elasticsearch	5
Transformation des Données pour l'Indexation	5
4. Processus d'Indexation	6
Préparation des Données	6
Indexation dans Elasticsearch	6
Indexation Batch et Unitaire	6
Automatisation de l'Indexation	7
Gestion des Erreurs	7
5. Conclusion	7
6. Annexes	8

1. Introduction

Contexte

Dans un monde de plus en plus connecté, la capacité de mettre en relation des individus ou des services en fonction de leur position géographique devient cruciale. Que ce soit pour des besoins sociaux, professionnels, ou de services de proximité, l'accès à des informations précises et en temps réel sur la localisation des individus ou des ressources offre des avantages significatifs. Notre projet s'inscrit dans cette dynamique, en proposant un outil de mise en relation par géolocalisation. Cet outil permettra aux utilisateurs de s'enregistrer, de fournir des informations personnelles ainsi que leurs coordonnées géographiques, et de rechercher des personnes disponibles dans un rayon défini.

Objectifs

Le projet vise à développer une application qui répond aux besoins suivants :

- Enregistrement des Utilisateurs : Permettre aux individus de s'enregistrer en fournissant des informations personnelles et leurs coordonnées de géolocalisation. L'enregistrement peut se faire soit individuellement, soit par lot à travers l'importation de fichiers JSON.
- Indexation des Données: Les données fournies par les utilisateurs seront indexées à l'aide d'un moteur d'indexation. Cette étape est cruciale pour optimiser les performances de recherche au sein de l'application.
- **Recherche géolocalisée :** Offrir une fonction de recherche permettant de trouver des utilisateurs en fonction de leur position géographique et dans un rayon spécifié par le demandeur.
- Affichage des Résultats: Les résultats de chaque recherche devront présenter les informations complètes des personnes trouvées, permettant ainsi des mises en relation efficaces et pertinentes dans un monde connecté, géolocalisation cruciale pour besoins sociaux/professionnels..

Technologies Utilisées

Pour répondre aux exigences fonctionnelles et techniques de ce projet, nous avons sélectionné un ensemble de technologies éprouvées et performantes :

- Langage de Développement : Python, choisi pour sa flexibilité, sa facilité d'utilisation et son vaste écosystème de bibliothèques, notamment pour le traitement des données et l'intégration avec les systèmes de bases de données et d'indexation.
- Moteur d'Indexation : Elasticsearch, retenu pour sa capacité à gérer de grandes quantités de données géolocalisées et à effectuer des recherches rapides et précises basées sur la géolocalisation.
- **SGBD**: MongoDB, une base de données NoSQL orientée documents, sélectionnée pour sa flexibilité dans le stockage et la gestion des données structurées de manière hiérarchique, ce qui est particulièrement adapté aux

- données personnelles et géolocalisées des utilisateurs.
- Système d'Exploitation: Linux, choisi pour sa stabilité, sa sécurité et sa performance, ainsi que pour sa large adoption dans les environnements de serveurs et de développement.

Le projet s'articule autour de ces technologies, chacune jouant un rôle clé dans la réalisation d'un système robuste, performant et capable de répondre aux besoins spécifiques de mise en relation géolocalisée.

2. Architecture Système

L'architecture est structurée de manière à faciliter les opérations d'enregistrement des utilisateurs, de stockage et d'indexation des données, ainsi que la recherche et l'affichage des informations. Les utilisateurs interagissent avec l'application principale via une interface utilisateur ou une API, l'application communique avec MongoDB pour le stockage des données et avec Elasticsearch pour l'indexation et la recherche géographique.

Voici un aperçu détaillé de chaque composant et de leur interaction.

Détail des Composants

- Application Principale (Python)
 - Fonctionnalités: Gère les interactions utilisateurs, l'enregistrement des données, les opérations de transformation et d'indexation, et les requêtes de recherche.
 - Intégration avec MongoDB: Enregistre et récupère les informations des utilisateurs, en utilisant le modèle de données conçu pour optimiser le stockage des informations complexes.
 - Intégration avec Elasticsearch : Transfère les données nécessaires depuis MongoDB vers Elasticsearch, applique les transformations requises pour l'indexation, et gère les requêtes de recherche géographique.
 - Interface Utilisateur/API: Offre aux utilisateurs un moyen d'interagir avec le système, que ce soit via une interface graphique ou une API pour l'enregistrement et la recherche.
- Base de Données (MongoDB)
 - Stockage des Données : Contient les informations détaillées des utilisateurs, y compris les données personnelles et géolocalisées, structurées selon le modèle de données spécifié.
 - Avantages : La flexibilité du modèle de données document et la capacité à gérer de grandes quantités de données avec des performances élevées.
- Moteur d'Indexation (Elasticsearch)
 - Indexation des Données : Indexe les informations des utilisateurs nécessaires pour la recherche géographique, en utilisant le mapping spécifié pour optimiser les requêtes.

- Recherche Géographique: Permet des recherches rapides et précises basées sur la localisation géographique, en utilisant des critères tels que la distance par rapport à un point donné.
- Système d'Exploitation (Linux)
 - Environnement d'Exécution : Fournit une plateforme stable et sécurisée pour l'exécution de l'application, de MongoDB et d'Elasticsearch, favorisant les performances et la fiabilité.

Interaction entre les Composants

Enregistrement des Utilisateurs : Les utilisateurs fournissent leurs informations via l'interface de l'application, qui les stocke ensuite dans MongoDB.

Transformation et Indexation : L'application extrait les données nécessaires de MongoDB, les transforme selon le modèle requis par Elasticsearch, et les indexe pour optimiser la recherche géographique.

Recherche : Lorsqu'une requête de recherche est soumise via l'interface utilisateur, l'application interroge Elasticsearch pour trouver les utilisateurs correspondants selon les critères géographiques. Les résultats sont ensuite récupérés de MongoDB pour afficher les informations complètes.

Affichage des Résultats : Les informations détaillées des utilisateurs trouvés sont présentées à l'utilisateur ayant initié la recherche.

Cette architecture offre une séparation claire des préoccupations entre le stockage des données, leur indexation pour la recherche géographique, et l'interaction utilisateur, assurant ainsi une gestion efficace des ressources et une expérience utilisateur optimale.

3. Modèle de Données

La conception du modèle de données est un élément fondamental de notre application de recherche géographique. Il détermine comment les informations des utilisateurs sont stockées, organisées et interconnectées. Pour ce projet, nous utilisons MongoDB comme base de données et Elasticsearch pour l'indexation des données géolocalisées. Voici le détail du modèle de données pour ces deux systèmes :

MongoDB

MongoDB, une base de données orientée documents, offre une grande flexibilité dans la modélisation des données. Les documents JSON (JavaScript Object Notation) permettent une structure hiérarchique, facilitant ainsi le stockage de données complexes comme celles des utilisateurs de notre application.

- Structure du Document Utilisateur :
 - id (long): Identifiant unique de l'utilisateur.
 - firstname (text): Prénom de l'utilisateur.

- lastname (text): Nom de famille de l'utilisateur.
- email (text): Adresse email de l'utilisateur.
- birthDate (date): Date de naissance de l'utilisateur.
- login (objet): Informations de connexion de l'utilisateur, comprenant un identifiant unique (uuid), un nom d'utilisateur (username), un mot de passe (password) et des informations de sécurité (md5, sha1, enregistrement).
- address (objet): Adresse de l'utilisateur, incluant la rue (street), le numéro d'appartement (suite), la ville (city), le code postal (zipcode) et les coordonnées géographiques (geo) avec latitude (lat) et longitude (lng).
- phone (text) : Numéro de téléphone de l'utilisateur.
- website (text): Site web de l'utilisateur.
- company (objet): Informations sur l'entreprise de l'utilisateur, incluant le nom (name), le slogan (catchPhrase) et le secteur d'activité (bs).

Ce modèle de document est conçu pour stocker de manière exhaustive les informations des utilisateurs, permettant ainsi une gestion complète et une recherche efficace au sein de l'application.

Cette étape de transformation garantit que les données sont optimisées pour la recherche géographique, tout en conservant toutes les informations nécessaires pour une mise en relation efficace entre les utilisateurs.

Elasticsearch

Elasticsearch utilise des structures de données appelées "indices" pour stocker les données de manière optimisée pour la recherche. Le mapping de notre indice users définit comment les documents doivent être stockés et indexés.

- Mapping de l'Indice users:
 - id (long): Identifiant unique de l'utilisateur, correspondant à l'id dans MongoDB.
 - name (text): Nom complet de l'utilisateur, combinant firstname et lastname de MongoDB.
 - username (text) : Nom d'utilisateur, extrait de l'objet login de MongoDB.
 - birthDate (date): Date de naissance de l'utilisateur.
 - address (text): Adresse complète de l'utilisateur, formatée à partir des données de l'objet address de MongoDB.
 - geo_point_2d (geo_point) : Coordonnées géographiques de l'utilisateur, stockées sous forme d'objet geo_point pour permettre des recherches géographiques efficaces dans Elasticsearch.

Le type <code>geo_point</code> est particulièrement important car il permet d'effectuer des recherches basées sur la localisation géographique, telles que trouver des utilisateurs dans un rayon spécifique autour d'un point donné.

Transformation des Données pour l'Indexation

La transformation des données de MongoDB vers Elasticsearch est cruciale pour assurer que les données soient correctement formatées et indexées pour la recherche. Ce processus implique :

- La combinaison des champs firstname et lastname de MongoDB pour créer le champ name dans Elasticsearch.
- La conversion de l'adresse de l'utilisateur en MongoDB en un format texte unique pour le champ address dans Elasticsearch.
- L'extraction et le formatage des coordonnées géographiques de MongoDB pour les adapter au type geo point dans Elasticsearch.

4. Processus d'Indexation

L'indexation des données dans notre système de recherche géographique est cruciale pour la performance et l'efficacité des recherches. Ce processus implique la transformation des données stockées dans MongoDB pour qu'elles soient optimisées pour la recherche dans Elasticsearch. Voici les étapes détaillées de ce processus d'indexation :

Préparation des Données

Avant l'indexation, les données doivent être préparées et transformées pour correspondre au mapping défini dans Elasticsearch. Cette préparation inclut :

- Extraction des Données de MongoDB : Sélection des documents utilisateurs à partir de MongoDB qui nécessitent une indexation ou une réindexation.
- Transformation des Données : Adaptation des données au format attendu par Elasticsearch. Cela inclut la combinaison des champs firstname et lastname en un seul champ name, la mise en forme de l'adresse et la conversion des coordonnées géographiques en un type geo point.

Indexation dans Elasticsearch

Une fois les données préparées, elles sont indexées dans Elasticsearch selon le processus suivant :

- Création/Mise à Jour de l'Index : Les documents transformés sont envoyés à Elasticsearch. Si un document avec l'ID spécifié existe déjà, il est mis à jour avec les nouvelles données; sinon, un nouveau document est créé.
- Utilisation du Type geo_point : Pour les champs de coordonnées géographiques, le type geo_point est utilisé, permettant à Elasticsearch d'optimiser les recherches basées sur la localisation.

Indexation Batch et Unitaire

Le système prend en charge deux modes d'indexation :

- Indexation Batch: Permet l'indexation de plusieurs documents en une seule opération. Typiquement utilisé lors de l'initialisation du système ou pour l'ajout en masse de nouveaux utilisateurs à partir de fichiers CSV ou JSON.
- Indexation Unitaire: Permet l'indexation d'un seul document à la fois, utilisé pour l'enregistrement ou la mise à jour des informations d'un utilisateur via l'interface utilisateur ou l'API.

Automatisation de l'Indexation

Pour maintenir l'index à jour avec les dernières données disponibles dans MongoDB, plusieurs stratégies peuvent être mises en œuvre :

- Indexation sur Modification : Chaque fois qu'un document utilisateur est ajouté ou modifié dans MongoDB, une opération d'indexation correspondante est automatiquement déclenchée pour mettre à jour l'index Elasticsearch.
- Indexation Périodique: Un processus planifié réindexe périodiquement les données dans Elasticsearch pour s'assurer que toutes les modifications dans MongoDB sont reflétées dans l'index.

Gestion des Erreurs

Durant l'indexation, le système doit être capable de gérer les erreurs potentielles, telles que les échecs de connexion à Elasticsearch ou les erreurs de format de données. Une stratégie robuste de gestion des erreurs inclut :

- Réessayer en Cas d'Échec : En cas d'échec temporaire, le système peut réessayer l'opération d'indexation un nombre défini de fois.
- Logging : Tous les échecs d'indexation sont enregistrés, permettant une analyse ultérieure pour résoudre les problèmes.

5. Conclusion

Le succès de ce projet repose sur une compréhension approfondie des besoins utilisateurs, qui a guidé la conception et le développement de chaque composant du système. En mettant l'accent sur une expérience utilisateur intuitive et des résultats de recherche précis et pertinents, nous avons créé un outil qui non seulement répond aux attentes des utilisateurs mais les dépasse car la structure du système, pensée pour être à la fois robuste et flexible, garantit non seulement la performance et la fiabilité de l'application mais offre également une plateforme évolutif pour intégrer des fonctionnalités supplémentaires à l'avenir.

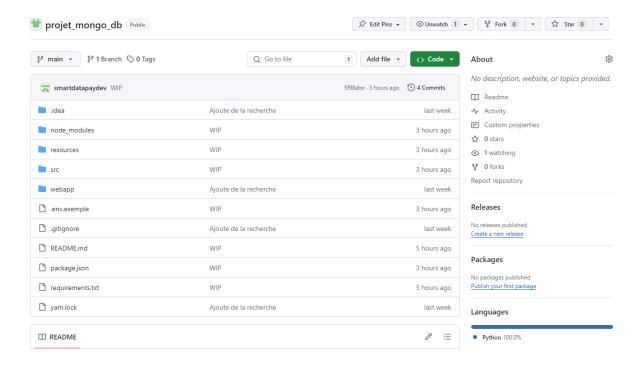
Ce projet démontre l'impact positif de l'ingénierie logicielle et des systèmes d'information modernes dans la résolution de problématiques complexes et la

création de valeur pour les utilisateurs. Il ouvre également la voie à de nouvelles explorations et innovations dans le domaine de la géolocalisation et de la mise en relation, offrant un vaste champ de possibilités pour les développements futurs. Notre engagement continu dans la recherche, le développement et l'amélioration de nos solutions est essentiel pour rester à l'avant-garde de la technologie et répondre aux défis émergents dans ce domaine dynamique.

6. Annexes

- Code source : Extraits significatifs du code développé.
- Captures d'écran : Illustrations de l'interface utilisateur et des exemples de résultats de recherche.

https://github.com/nkaurelien-ionis-stm/projet_mongo_db



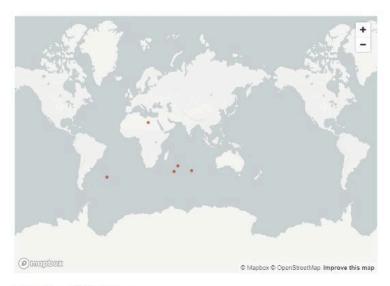
Recherche Géographique dans Elasticsearch



Resulat de la recherche:

Nombre d'utilisateurs trouvées sur 10000km: 5 / 5

	Nom	Username	Adresss	Latitude	Longitude
0	Leanne Graham	Bret	Kulas Light, 92998-3874 Gwenborough	-37.3159	81.1496
1	Ervin Howell	Antonette	Victor Plains, 90566-7771 Wisokyburgh	-43,9509	-34.4618
2	Chelsey Dietrich	Kamren	Skiles Walks, 33263 Roscoeview	-31.8129	62.5342
3	Kurtis Weissnat	Elwyn.Skiles	Rex Trail, 58804-1099 Howemouth	24.8918	21.8984
4	Clementina DuBuque	Moriah.Stanton	Kattie Turnpike, 31428-2261 Lebsackbury	-38.2386	57.2232



Tous les utilisateurs

10 / 10

Clementina DuBuque	Moriah.Stanton	Kattie Turnpike, 31428-2261	-38.2386
Glenna Reichert	Delphine	Dayna Park, 76495-3109 Bart	24.6463
Nicholas Runolfsdottir V	Maxime_Nienow	Ellsworth Summit, 45169 Ali	-14.3990
Kurtis Weissnat	Elwyn.Skiles	Rex Trail, 58804-1099 Howe	24.8918
Mrs. Dennis Schu <mark>l</mark> ist	Leopoldo_Corkery	Norberto Crossing, 23505-13	-71.4197
Chelsey Dietrich	Kamren	Skiles Walks, 33263 Roscoevi	-31.8129
Patricia Lebsack	Karianne	Hoeger Mall, 53919-4257 So	29.4572
Clementine Bauch	Samantha	Douglas Extension, 59590-41	-68.6102
Ervin Howell	Antonette	Victor Plains, 90566-7771 Wi	-43.9509
Leanne Graham	Bret.	Kulas Light, 92998-3874 Gwe	-37.3159
Nom	Username	Adresss	Latitude

```
{ □
      "_index":"users",
      "_id": "CYgxm44BYzzxx51UPpQ9",
      "_score":1.0,
      "_source":{ ⊟
         "id":1,
         "name": "John Doe",
         "username": "johndoe",
         "birthDate": "1973-01-22",
         "address": "123 Main Street, 12345-6789 Anytown\nApt. 4",
         "geo_point_2d":{
            "lat": "42.1234",
            "lon": "-71.2345"
   },
   { □
      "_index": "users",
      "_id": "Cogxm44BYzzxx51UPpQ9",
      "_score":1.0,
      "_source":{ ⊟
         "id":2.
         "name": "Jane Smith".
         "username": "janesmith",
         "birthDate": "1983-02-22",
         "address": "456 Oak Street, 12345-6789 Anytown\nSuite 200",
         "geo_point_2d":{
            "lat": "42.3456",
            "lon":"-71.6789"
         }
     }
   },
```

```
NosqlMongoDb > resources > data > {} users.json > ...
       You, 1 hour ago | 1 author (You)
  1
       "id": 1,
           "firstname": "John",
           "lastname": "Doe",
           "email": "johndoe@example.com",
           "birthDate": "1973-01-22",
           "login": {
             "uuid": "1a0eed01-9430-4d68-901f-c0d4c1c3bf22",
             "username": "johndoe",
             "password": "jsonplaceholder.org",
             "md5": "c1328472c5794a25723600f71c1b4586",
             "sha1": "35544a31cc19bd6520af116554873167117f4d94",
             "registered": "2023-01-10T10:03:20.022Z"
           "address": {
             "street": "123 Main Street",
             "suite": "Apt. 4",
             "city": "Anytown",
             "zipcode": "12345-6789",
             "geo": {
               "lat": "42.1234",
               "lng": "-71.2345"
           "phone": "(555) 555-1234",
           "website": "www.johndoe.com",
           "company": {
             "name": "ABC Company",
             "catchPhrase": "Innovative solutions for all your needs",
             "bs": "Marketing"
```

```
NosqlMongoDb > $ .env.exemple
You, 3 hours ago | 1 author (You)

ES_SERVER_NAME=3.72.35.55

ES_SERVER_PORT=9200

MONGO_DB_DATABASE=projet_mongo_db

MONGO_DB_URL = "mongodb+srv://${MONGO_DB_USER}:${MONGO_DB_PASSWORD}@cluster0.xifg

| Comparison of the comparison of
```

```
NosqlMongoDb > src > ♥ streamlitappuser.py > ...
      import streamlit as st
      import pandas as pd
      from st aggrid import AgGrid
      from utils.es import Search
      es = Search()
      es.index_name = "users"
      columns = ["Nom", "Username", "Adresss" , "Latitude", "Longitude" ]
      start_lat=24.8918
      start_lng=21.8984
 14
      @st.cache_data
      def get_data() -> pd.DataFrame:
           (users, total) = es.get_all_data()
           print(users)
          data = []
          df = pd.DataFrame(data, columns=columns)
           if users:
               for hit in users:
                   source = hit["_source"]
                   data.append([source["name"], source["username"], source["address"], source["geo_point_
                                 source["geo_point_2d"]["lon"]])
               if data:
                   df = pd.DataFrame(data, columns=columns)
           return (df, total)
      st.title('Recherche Géographique')
      # Formulaire de recherche géographique
      with st.form(key='search_geo'):
          lat = st.number_input('Latitude', value=start_lat) # Valeur par défaut pour Kurtis Elwyn
          lon = st.number_input('Longitude', value=start_lng) # Valeur par défaut pour Kurtis Elwyn
rayon = st.number_input('Rayon (km)', value=10000)
           submit_button = st.form_submit_button(label='Rechercher')
```

```
NosqlMongoDb > src > utils > ♦ mapping.py > ...
        def transform_data(document, index_name):
             transformed_document = document # This is just a placeholder line
             if index_name == 'gares' :
                 transformed_document = {
                      "id": document["code_uic"],
                       "libelle": document["libelle"],
                      "commune": document["commune"],
"departemen": document["departemen"],
                       "geo_point_2d": {"lat": document["geo_point_2d"]["lat"], "lon": document["geo_point_2d"]["lon"]}
             elif index_name == 'users' :
                   transformed_document = {
                       "id": document["id"],
                       "name": "{} {}".format(document["firstname"], document["lastname"]),
                       "username": document["login"]["username"],
                      "birthDate": document["birthDate"],

"address": "{}, {} {}\n{}".format(document["address"]["street"], document["address"]["zipcode"], doc

"geo_point_2d": {"lat": document["address"]["geo"]["lat"], "lon": document["address"]["geo"]["lng"]}
             return transformed_document
 27
```

```
NosqlMongoDb > src > utils > ♦ import_mongodb.py > ...
      from pymongo.mongo_client import MongoClient
      from pymongo.server_api import ServerApi
      import os
      import json
      from dotenv import load_dotenv
      load_dotenv()
      db_name=os.environ["MONGO_DB_DATABASE"]
      mongo_db_server_url = os.environ["MONGO_DB_URL"]
      # Create a new client and connect to the server
      client = MongoClient(mongo_db_server_url, server_api=ServerApi('1'))
      db = client[db_name]
      directory = 'resources/data' # Update this path to the directory containing your
      # List all JSON files in the directory
      files = [f for f in os.listdir(directory) if f.endswith('.json')]
      for file in files:
          collection_name = os.path.splitext(file)[0]
          collection = db[collection_name]
          # Drop the collection if it exists
          collection.drop()
          with open(os.path.join(directory, file), 'r') as json_file:
              data = json.load(json_file)
              # Assuming each file contains an array of documents
              if isinstance(data, list):
                  collection.insert_many(data)
              else: # If the file contains a single document
                  collection.insert_one(data)
```

```
NosqlMongoDb > src > utils > ♦ create_search_index.py > ...
 27
      # Import each JSON file into MongoDB
      for data_file in data_files:
          # Extracting the filename without extension to use as collection name
           index name = os.path.splitext(data file)[0]
          if es.indices.exists(index=index_name):
               # Delete the index if it exists
               es.indices.delete(index=index_name)
          mapping file_path = os.path.join(mapping_directory, index_name + '.json')
          if os.path.exists(mapping file path):
              with open(mapping file path, 'r') as mapping file:
                  mapping = json.load(mapping_file)
                   # es.indices.create(index=index_name, body=mapping, ignore=400) # ignore
                  es.indices.create(index=index_name, body=mapping)
               print(f"No mapping file found for index {index_name}. Creating index without
               # Create the index without specific mapping
              es.indices.create(index=index_name)
          with open(os.path.join(data_directory, data_file), 'r') as f:
               data = json.load(f)
               # Assume data is a list of documents. If not, adjust accordingly.
              transformed_data = [transform_data(doc, index_name) for doc in data]
               actions = [
                       " index": index_name,
                       " source": doc,
                   for doc in transformed data
              helpers.bulk(es, actions)
      print("Data transformation and insertion complete.")
```

```
NosqlMongoDb > (i) README.md > • # INSTALL
      # INSTALL
      Install python virtual env
      ```console
 pip install --upgrade pip
 sudo apt install python3-virtualenv
 pip install --upgrade virtualenv
 11
 # Create virtual Env
 Install python lib
      ```shell
      virtualenv -p python3 venv
      chmod +x venv/bin/activate
      Install python lib
      ```shell
 pip install -r requirements.txt
 # RUN
 24
 Activate python venv
      ```shell
      source ./venv/bin/activate
 29
      ### Import dataset to mongoDb
      ```shell
 python src/utils/import_mongodb.py
 ### Import dataset to Elasticseacrh
      ```shell
      python src/utils/create_search_index.py
      # Run app
      ```shell
 41
 yarn start
 streamlit run src/streamlitappuser.py
```