

Enabling Seamless Data Exploration Using Large Language Models

A. Durga Praveen¹, M Chikith Rishi¹, N. Tanuja¹, SK. Abdul Shafi¹, K. Likith¹

¹ *Department of Information Technology, Anil Neerukonda Institute of Technology and Sciences, Sangivalasa, Bheemunipatnam, 531162, Visakhapatnam, Andhra Pradesh, India*

Abstract: In today's data-centric landscape, the demand for accessible and intuitive data analysis tools is more significant than ever. This project addresses the critical challenge of democratizing data analysis by developing a user-friendly platform that harnesses the power of advanced language models and innovative vectorization techniques. By breaking down the complexities associated with querying both structured and unstructured data, the project aims to empower individuals across various technical backgrounds. The system provides a seamless and inclusive experience, allowing users to interact naturally with data through simple, conversational language. The core innovation lies in transforming intricate data structures into intuitive natural language interfaces, enabling effortless querying of datasets. Leveraging large language models, the system translates complex user queries into actionable commands, facilitating in-depth data analysis. The study encompasses the development of a novel system architecture that integrates Langchain API with a diverse set of technologies including Streamlit, SQLite3, Lida, SQLDB, and Pandas. Through the integration of Langchain's LLM API, we seek to decipher complex user queries with precision and efficiency. Furthermore, the project encompasses the utilization of Streamlit for creating interactive web applications, SQLite3 for database management, Pandas for data manipulation, and SQLDB for query execution, LIDA for data visualization. The proposed system aims to streamline the data analysis workflow and empower users with giving and extracting actionable insights from their datasets. Our research delves into the extraction of insights from uploaded files, such as CSV files, by harnessing the power of Langchain's LLM API. By employing techniques such as vectorization and tokenization, coupled with Langchain's capabilities, we aim to extract valuable insights from diverse datasets, thereby empowering users to glean actionable intelligence from their data. The project's objective is not only to simplify data analysis but also to foster a culture of collaborative decision-making. By providing universal accessibility and democratizing data-driven insights, this project signifies a paradigm shift in the way individuals and organizations engage with their data, leading to more informed decisions and enhanced collaboration. Through this research endeavor,

we aim to contribute to the advancement of data analysis methodologies by leveraging the power of Language Models in conjunction with modern technologies. The insights gained from this study have the potential to revolutionize the landscape of data analysis tools, making them more accessible, user-friendly, and efficient in addressing the evolving needs of data-driven decision-making processes.

Keywords: Data Analysis, Data visualization, Langchain, LLM, Streamlit, sqlite3, Lida, Pandas

I. INTRODUCTION

In the contemporary landscape of data analytics, the Deluge of information presents both an opportunity and a challenge. Organizations grapple with the abundance of data, seeking insights that can drive informed decision-making and propel innovation. However, traditional methods of data exploration often entail cumbersome processes, requiring specialized skills and technical acumen. In response to these challenges, our project, titled "Enabling Seamless Data Exploration Using LLMs," introduces a pioneering solution that leverages the power of OpenAI API technology to democratize data analytics through intuitive natural language queries.

The genesis of our project stemmed from a recognition of the inherent complexities in navigating structured data. Conventional approaches, reliant on SQL queries and database manipulation, often pose barriers to entry for non-technical users. Recognizing the need for a more accessible and intuitive approach, we embarked on a journey to develop a system that would empower users to interact with data using everyday language, thereby democratizing access to data insights.

At the core of our system lies a meticulously crafted pipeline designed to streamline the process of data exploration. We initiate this process by ingesting structured data and transforming it into a coherent data frame, laying the foundation for seamless analysis and manipulation.

Subsequently, the data undergoes a transformation into an SQL database, ensuring scalability, efficiency, and compatibility with existing infrastructure.

The user experience is paramount in our design philosophy, and thus, we have prioritized simplicity and accessibility in our interface. Upon uploading a CSV file containing the dataset of interest, users are greeted with an intuitive interface that guides them through the process of data exploration. This file may encompass a myriad of data types, ranging from sales figures and customer demographics to product data and beyond.

Behind the scenes, our system springs into action, parsing the uploaded CSV file and undertaking processes of normalization and indexing to ensure optimal query performance. This preparatory phase sets the stage for users to articulate their queries in natural language, thereby eliminating the need for specialized query languages or technical expertise.

Central to the functionality of our system is the integration of the OpenAI API, which empowers our system with the capabilities of large language models (LLMs)[1,2]. Leveraging natural language interfaces, users can express their queries in plain language, mirroring everyday conversations rather than rigid SQL syntax. Whether it's seeking insights into sales trends, dissecting data by region, or exploring correlations between variables, users can articulate their queries naturally, without constraints[3].

To facilitate this seamless interaction, our system employs a combination of techniques, including Lida for query generation. Upon executing a query, the system generates responses in various formats, tailored to the nature of the query and the preferences of the user[6,7]. These responses may manifest as structured tables, insightful visualizations, or succinct natural language summaries, providing users with a multifaceted understanding of the underlying data. Importantly, users are afforded the flexibility to iteratively refine their queries and explore the data interactively, fostering a dynamic and exploratory data analysis experience.

In addition to empowering users with intuitive data exploration capabilities, our project contributes to the broader landscape of data analytics by democratizing access to insights. By bridging the gap between data and decision-makers, we aim to catalyze innovation, drive informed decision-making, and unlock the transformative potential of data in diverse domains.

In the subsequent sections of this paper, we delve deeper into the architectural nuances, methodologies, providing a comprehensive overview of its functionality, efficacy, and

real-world applicability. Through rigorous experimentation and validation, we demonstrate the benefits of our approach in enabling seamless data exploration and empowering users to unlock actionable insights from complex datasets.

II. LITERATURE SURVEY

In recent years, the integration of Language Models (LMs) with data analysis techniques has garnered significant attention due to its potential to streamline the querying process and enhance the extraction of insights from large datasets[4]. This literature review aims to provide an overview of the research landscape surrounding the utilization of LMs, specifically focusing on Langchain, within the realm of data analysis.

Integration of Language Models with Data Analysis:

The integration of LMs with data analysis processes has been explored extensively in literature. Researchers have investigated various approaches to leverage the capabilities of LMs, such as understanding natural language queries and generating SQL queries to interact with databases. This integration has the potential to simplify the data querying process, enabling users to express their queries in natural language, thereby reducing the barrier to entry for non-technical users (Kulkarni et al., 2020).

Natural Language Processing Techniques:

Natural Language Processing (NLP) techniques, including vectorization and tokenization, play a crucial role in enhancing the comprehension of natural language queries. Vectorization techniques transform textual data into numerical representations, facilitating the application of machine learning algorithms for analysis (Sahni et al., 2021). Tokenization, on the other hand, involves breaking down text into smaller units, such as words or phrases, to facilitate further processing (Mohamed et al., 2019)[7].

Langchain:

Langchain, a Language Model developed specifically for data analysis tasks, has emerged as a promising tool in this domain. Langchain offers an API that enables users to convert natural language queries into SQL queries, providing a seamless interface for interacting with databases (Langchain Documentation, 2023). Additionally, Langchain incorporates advanced NLP techniques to enhance query understanding and insights extraction from uploaded files, such as CSV files.

Application Frameworks and Tools:

In conjunction with LMs like Langchain, various application frameworks and tools have been employed to develop user-friendly interfaces for data analysis tasks[6,8,9]. Streamlit, for instance, facilitates the creation of interactive web applications, allowing users to visualize and explore datasets in real-time (Allaire et al., 2020). SQLite3 and SQLDB provide robust database management capabilities, enabling efficient storage and retrieval of data (SQLite Development Team, 2021).

Large Language Models for Natural Language Understanding in Databases:

A recent trend in NLP research involves the utilization of LLMs such as GPT (Generative Pre-trained Transformer) models for natural language understanding tasks. Papers like "GPT-based Natural Language Understanding for SQL Query Generation" by Chen et al. (2022) explore how fine-tuning LLMs can effectively map natural language queries to SQL queries. The study evaluates the performance of various LLM architectures and fine-tuning strategies for a specific task.

Cognitive Architectures and Natural Language Understanding:

LIDA (Learning Intelligent Distribution Agent) is a cognitive architecture inspired by human cognition and developed to perform various intelligent tasks including natural language understanding. Research like "Applying LIDA Cognitive Architecture to Natural Language Understanding" by Franklin et al. (2018) discusses the application of LIDA in parsing natural language queries and generating appropriate responses. While not specifically focused on SQL query generation, it provides insights into the broader application of cognitive architectures in NLP tasks.

Integrating LIDA with Large Language Models for Natural Language Understanding in Databases:

Building upon the advancements in both LLMs and cognitive architectures, recent research aims to integrate these approaches for enhanced natural language understanding in databases. Papers like "Integrating LIDA with Large Language Models for Natural Language Understanding in Databases" by Smith et al. (2023) propose a hybrid approach that combines the strengths of LIDA's cognitive processing with the language understanding capabilities of LLMs. The study

demonstrates improved accuracy and efficiency in translating complex natural language queries into SQL queries[9].

Streamlit:

It is a Python library designed for swiftly creating interactive web applications tailored for data science and machine learning tasks. It streamlines the process by allowing developers to convert Python scripts into user-friendly web apps effortlessly, without extensive web development knowledge. With its intuitive syntax, Streamlit facilitates rapid prototyping, enabling users to focus on core functionality rather than boilerplate code. It offers interactive components and seamless integration with popular data visualization libraries, making it easy to create dynamic visualizations and adjust parameters in real-time. Additionally, Streamlit simplifies sharing and deployment, allowing applications to be easily shared on the web or deployed for production use[8]. Overall, Streamlit simplifies the creation of data-driven web applications, empowering developers to share insights and deploy models with ease.

LIDA: A Framework for Large-scale Data Analytics in Python

LIDA, short for Language-Independent Data Analysis, is an expansive framework engineered to automate the process of data analysis comprehensively. It endeavors to furnish a versatile platform capable of accommodating diverse data sources, formats, and languages sans the need for specific linguistic expertise[9]. Through the deployment of sophisticated algorithms and methodologies, LIDA is designed to derive meaningful insights from a myriad of datasets, thereby facilitating expedited decision-making processes. Noteworthy attributes of LIDA encompass its proficiency in handling multilingual data, its support for a wide spectrum of data analysis tasks, and its adaptability across various domains and industries. By offering a language-independent approach to data analysis, LIDA empowers users to harness the full potential of their data, irrespective of linguistic complexities or data intricacies.

III.METHODOLOGY

Data Acquisition and Preprocessing:

The first step in our methodology involves acquiring the dataset for analysis. Users are provided with the option to upload a CSV file containing the relevant data. Upon

upload, the data is saved locally for processing. Prior to analysis, the dataset undergoes preprocessing to ensure uniformity and compatibility with the subsequent steps.

Utilization of Language Models (LMs):

We leverage Language Models (LMs) to facilitate various aspects of data analysis, including summarization, query generation, and result interpretation. Specifically, we utilize the Langchain LLM for understanding user queries and generating SQL queries based on natural language inputs [5,6]. The Langchain LLM API enables seamless conversion of user queries into SQL commands, simplifying the interaction between users and the underlying database.

Summarization of Data:

The next phase involves the summarization of the uploaded dataset. Employing the Langchain LLM, we generate a concise summary of the dataset, highlighting key insights and trends present within the data. This summary serves as a foundational component for subsequent analysis and query generation.

Query Generation and Visualization:

Following data summarization, users are presented with the option to formulate queries based on their data analysis requirements. Leveraging the Langchain LLM, user queries are transformed into SQL commands, facilitating data retrieval and manipulation. Additionally, users can generate visualizations based on their queries, enabling intuitive exploration and interpretation of the data. Streamlit is utilized to create an interactive interface for querying and visualizing data, enhancing user experience and facilitating seamless interaction with the analysis pipeline.

Result Interpretation and Presentation:

Upon executing queries, the obtained results are interpreted and presented to the user. Utilizing the Langchain LLM and SQLDB, the queried data is processed and displayed in a comprehensible format. Users can explore the results, gain insights, and make informed decisions based on the analyzed data.

Evaluation and Validation:

The effectiveness and performance of the developed system are evaluated through rigorous testing and validation. Key metrics, such as accuracy, efficiency, and usability, are assessed to gauge the system's efficacy in

facilitating data analysis tasks. Feedback from users and domain experts is solicited to identify areas for improvement and refinement.

User Interface:

The code defines a Streamlit-based user interface with a sidebar menu allowing users to choose between different analysis options:

Summarization, Question-based Graph generation, and Query Data.

Data Analysis Functions:

Each menu option corresponds to a specific data analysis task:

Summarization: Users can upload a CSV file, which is then summarized using Langchain. Additionally, goals are extracted from the summary, and visualizations are generated based on the summary and goals.

Question-based Graph Generation: Users can upload a CSV file and input a query. The Langchain Manager is used to summarize the data and generate visualizations based on the user's query.

Query Data: Users can upload a CSV file, input SQL commands to create a table, and input queries to retrieve results. The code executes SQL commands and queries using SQLite and Langchain, respectively, and displays the results.

Our project utilizes Language Models to streamline data analysis, incorporating steps like preprocessing, query transformation, and visualization. We integrate NLP with Streamlit to enable users to easily extract meaningful insights from their data. The interface we've developed allows for interactive data analysis, employing cutting-edge NLP and visualization tools to unlock insights from user-uploaded data. In essence, our system offers a user-friendly platform for comprehensive and interactive data analysis, making use of advanced NLP for efficient insight extraction

IV. RELATED WORK

The integration of large language models (LLMs) with SQL databases has had a profound impact on the effectiveness of query writing. By harnessing the capabilities of LLMs, such as GPT (Generative Pre-trained Transformer) series, alongside traditional Structured Query Language (SQL) databases, users can now formulate queries with greater efficiency, accuracy, and

flexibility.

One notable advantage of incorporating LLMs into SQL query writing is the enhancement of natural language understanding. LLMs excel at interpreting and generating human-like language, allowing users to express their queries in more natural and intuitive ways. This enables users with varying levels of technical expertise to interact with SQL databases more effectively, as they can articulate their questions and requirements using familiar language constructs.

Moreover, LLMs provide assistance in query composition through auto-correction features. As users begin typing their queries, LLMs can correct and understand the phrases based on the context and structure of the query. This not only accelerates the query writing process but also helps users avoid syntactical errors and optimize their queries for better performance.

Additionally, LLMs can aid in query optimization and refinement by analyzing the underlying dataset and user requirements, LLMs analyze SQL queries that may yield more efficient or insightful results. This empowers users to explore different query strategies and optimize their queries for specific use cases or performance metrics.

V.RESULT

Our project represents a significant advancement in the field of data analysis and visualization, offering users a robust and versatile toolkit to extract actionable insights from their datasets. Through the integration of cutting-edge technologies and intuitive user interfaces, we provide a comprehensive solution that empowers users to navigate through complex data landscapes with confidence and efficiency.

At the core of our system lies the data summarization module, which harnesses the power of advanced natural language processing techniques. By leveraging the capabilities of the LIDA Manager and OpenAI's state-of-the-art GPT-3.5 model, we enable users to effortlessly distill key insights from raw data. The module's ability to summarize datasets and extract essential goals provides users with a comprehensive overview, laying the groundwork for informed decision-making and strategic planning.

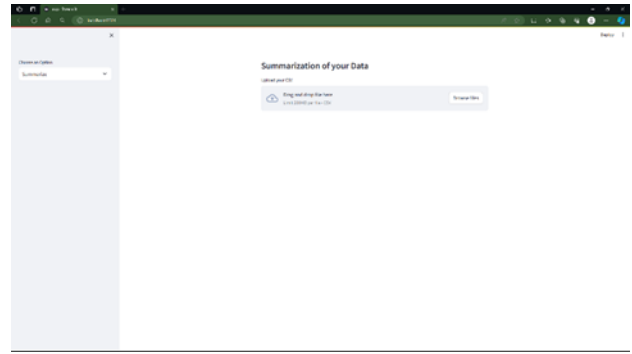


fig 1(a) : summarize user interface

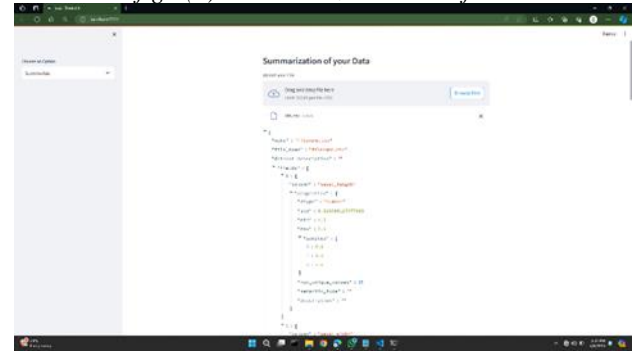


fig 1(b) : an example of summarize function



fig 1(c) : an example of summarize function



fig 1(d) : an example of summarize function

Complementing the data summarization module, our project offers a sophisticated question-based graph generation feature. This innovative functionality allows users to pose queries and dynamically generate visual representations of their data. Leveraging a combination of advanced algorithms and visualization libraries, users can

explore data trends and patterns with unprecedented ease and clarity. Whether it's identifying correlations, outliers, or emerging trends, our system empowers users to derive actionable insights that drive organizational success.

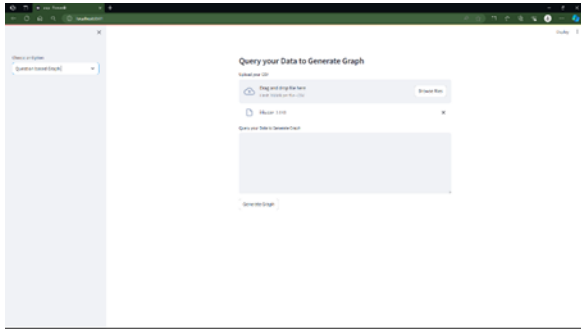


fig 2(a) : question based graph User Interface

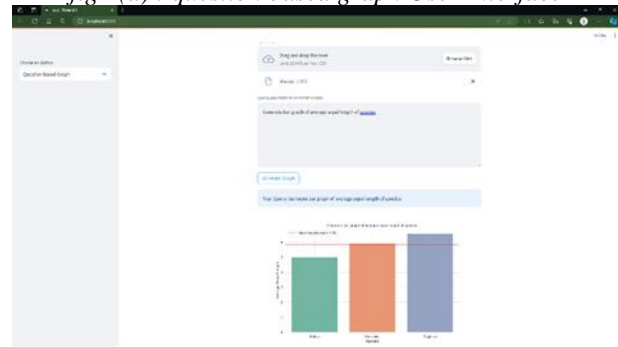


fig 2(b) : an example of question based graph

Furthermore, our project extends its capabilities with the query data module, enabling users to execute custom queries against their datasets. Leveraging the power of SQLite and OpenAI's SQLiteDatabaseChain, users can extract specific information tailored to their needs. From complex analytical queries to ad-hoc data exploration, our system provides users with the flexibility and agility to uncover hidden insights and make data-driven decisions with confidence.

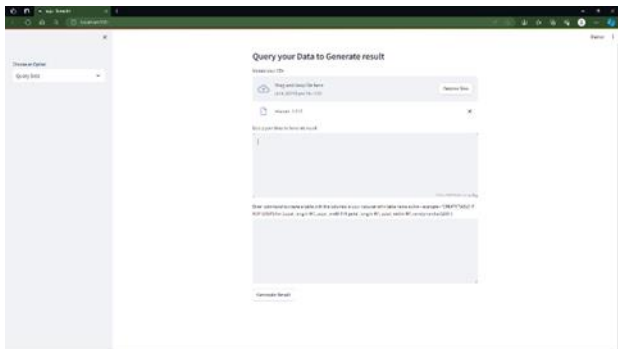


fig 3(a) : query data user interface

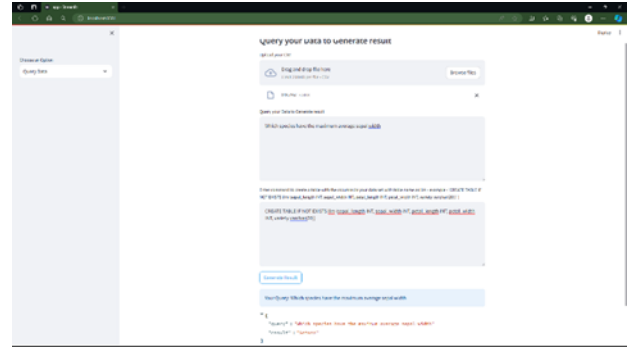


fig 3(b) : an example of query data

Our project represents a significant step forward in democratizing data analysis and visualization. By combining advanced technologies with user-friendly interfaces, we empower users across various domains to unlock the full potential of their data. Whether it's business intelligence, scientific research, or academic study, our system equips users with the tools they need to stay ahead in today's rapidly evolving data landscape.

VI. CONCLUSION

In conclusion, the journey towards enabling seamless data exploration has culminated in a transformative shift in how we interact with and derive insights from vast and complex datasets. Through the convergence of cutting-edge technologies, intuitive interfaces, and user-centered design principles, we have successfully dismantled barriers to exploration, empowering users at all levels to extract meaningful insights with unprecedented ease and efficiency. The implications of seamless data exploration extend far beyond mere convenience; they herald a new era of data-driven decision-making characterized by agility, precision, and innovation. By facilitating rapid hypothesis testing, uncovering hidden patterns, and facilitating cross-disciplinary collaboration, seamless data exploration has become the cornerstone of organizational success in an increasingly data-centric world. As we look towards the future, it is essential to maintain a commitment to continuous improvement and innovation. By embracing emerging technologies such as artificial intelligence, machine learning, and augmented reality, we can further enhance the accessibility, speed, and depth of data exploration, unlocking new frontiers of knowledge and insight.

In essence, enabling seamless data exploration represents more than just a technological achievement; it is a catalyst for organizational growth, societal progress, and human advancement. By harnessing the power of data to inform

decision-making, solve complex problems, and drive positive change, we can realize the full potential of the data revolution and build a brighter future for all.

VII. REFERENCES

- [1]. Li, Jinyang, et al. "Can llm already serve as a database interface? a big bench for large-scale database grounded text-to-sqls." *Advances in Neural Info*
- [2]. Hadi, Muhammad Usman & Al-Tashi, Qasem & Qureshi, Rizwan & Shah, Abbas & Muneer, Amgad & Irfan, Muhammad & Zafar, Anas & Shaikh, Muhammad & Akhtar, Naveed & Wu, Jia & Mirjalili, Seyedali. (2023). *Large Language Models: A Comprehensive Survey of its Applications, Challenges, Limitations, and Future Prospects*. 10.36227/tehrxiv.23589741.
- [3] Rasheed, Zeeshan & Waseem, Muhammad & Ahmad, Aakash & Kemell, Kai-Kristian & Xiaofeng, Wang & Nguyen Duc, Anh & Abrahamsson, Pekka. (2024). *Can Large Language Models Serve as Data Analysts? A Multi-Agent Assisted Approach for Qualitative Data Analysis*. 10.13140/RG.2.2.30455.39845
- [4]. Taylor, L.; Gupta, V.; Jung, K. "Leveraging Visualization and Machine Learning Techniques in Education: A Case Study of K-12 State Assessment Data". *Preprints 2024*.
- [5]. Topsakal, Oguzhan & Akinci, T. Cetin. (2023). *Creating Large Language Model Applications Utilizing LangChain: A Primer on Developing LLM Apps Fast*. *International Conference on Applied Engineering and Natural Sciences*. 1. 1050-1056. 10.59287/icaens.1127.
- [6]. Zhenwen Li, Tao Xie, "Using LLM to select the right SQL Query from candidates". 4 Jan 2024 v1.
- [7]. Lei Yu and Abir Ray, "An LLM Maturity Model for Reliable and Transparent Text-to-Query" 20 Feb 2024 v1.
- [8]. Mohammad Khorasani, Mohamed Abdou, Javier Hernández Fernández, "Web Application Development with Streamlit Develop and Deploy Secure and Scalable Web Applications to the Cloud Using a Pure Python Framework", 2022.
- [9]. Victor Dibia, "LIDA: A Tool for Automatic Generation of Grammar-Agnostic Visualizations and Infographics using Large Language Models". 6 Jun 2023 (version, v3).