



GenInsight

Project Proposal

Project Advisor:

Ms. Sana Fatima

Group Members:

Alishba Kamran 21L-6297

Anas Tanvir 21L-5678

Usman Bin Imran 21L-7963

National University Of Computer and Emerging Sciences
Department of Computer Science
Lahore, Pakistan

Abstract

GenInsight is a comprehensive, cloud-hosted data analytics platform designed to revolutionize data management, advanced analytics, and user interaction in the era of Industry 5.0, where data is the new oil. Built on a robust cloud infrastructure, GenInsight ensures high availability, load balancing, and efficient data handling. The platform integrates a secure user management system with state-of-the-art cybersecurity measures, including password hashing and privacy protection, guaranteeing the safety of user data.

Leveraging cache memory, GenInsight delivers fast and responsive user experiences. It features an advanced data analytics tool for in-depth Exploratory Data Analysis (EDA), as well as Adhoc Analysis and Adhoc EDA, enabling users to explore and manipulate data dynamically. The platform's customizable dashboard, powered by Natural Language Processing (NLP) and Large Language Model (LLM) technologies, allows for intuitive, conversational interactions with data.

GenInsight also includes a dedicated chat page with voice integration and an automated machine learning (AutoML) module that intelligently selects the optimal models for the data. This plug-and-play tool is essential for organizations aiming to harness the full potential of their data, streamlining preprocessing, analysis, and visualization, and empowering data-driven decision-making in the fast-paced world of Industry 5.0

1. Introduction

GenInsight is a cloud-based data management and analytics platform being developed as our Final Year Project. Addressing the shift towards cloud-hosted data, this project leverages cloud computing's scalability, load balancing, and security to deliver a robust solution. Designed to be user-friendly for all company members, including C-suite executives, board members, managers, and the chairman, GenInsight ensures that even non-technical users can easily navigate and utilize its features. The platform offers advanced tools for data visualization and analysis, including Adhoc Exploratory Data Analysis (EDA) and Adhoc reporting capabilities. It incorporates technologies like cache memory for speed optimization and Natural Language Processing (NLP) for customizable dashboards. Additionally, voice-integrated chat enhances real-time collaboration, making data interaction seamless and intuitive for everyone in the organization.

2. Goals and Objectives

The primary objectives of our project are:

- To create a cloud-based platform for efficient data management and analytics.
- To implement a secure user management system with advanced cybersecurity measures.
- To optimize application performance using cache memory.
- To provide an advanced EDA tool and customizable dashboard using NLP.
- To enable real-time collaboration through a chat feature with voice integration.
- To facilitate automated machine learning for selecting the best models for the data.
- To streamline data preprocessing for seamless integration with analytics tools.

3. Scope of the Project

GenInsight aims to create a cloud-based platform for data management, analysis, and visualization with the following components:

- A user-friendly interface with HTML, CSS, and JavaScript for smooth navigation and data interaction.
- Python for core logic, data processing, and backend communication to support analytics and machine learning.
- Azure Cloud for scalable, secure hosting and efficient data management.
- Tools for in-depth data analysis and reporting using Python libraries.
- NLP for user interaction through natural language queries.
- Real-time chat with voice integration for seamless user collaboration.
- AutoML for automated model selection and insight generation.
- Robust data cleaning, transformation, and preparation features.

4. Initial Study and Work Done so Far

Research into cloud-based data analytics platforms has highlighted advancements in data security, scalability, NLP, and AutoML, which are crucial for developing GenInsight.

Armbrust et al. [1] emphasize the benefits of cloud computing, such as scalability and high availability. GenInsight leverages Microsoft Azure's infrastructure for its data management and scalability needs, ensuring efficient handling of large datasets.

Stallings [2] underscores the importance of secure authentication and encryption. GenInsight implements OAuth 2.0 and bcrypt to protect user data and ensure secure access.

Cattell [3] discusses caching mechanisms to enhance performance. GenInsight uses Redis for caching, reducing latency and speeding up data retrieval, which is essential for handling complex datasets.

Bostock et al. [4] and Bird et al. [5] provide foundational technologies for data visualization and NLP. GenInsight employs pandas, matplotlib, and seaborn for EDA, and integrates spaCy and GPT-4 for NLP, allowing users to interact with data through natural language queries.

AutoML principles from Hutter et al. [6] guide GenInsight's model selection and optimization. Azure AutoML automates the process, improving accuracy and efficiency in data analysis.

Prototyping has focused on secure authentication, data caching, EDA, and NLP-based dashboard customization. These prototypes have validated core functionalities and highlighted areas for further development, such as scaling strategies and user interface enhancements.

These efforts form a strong foundation for GenInsight, with ongoing research and testing guiding its evolution to meet user needs and industry standards.

References

- [1] M. Armbrust et al., "A View of Cloud Computing," *Communications of the ACM*, 2010. Available at: <https://dl.acm.org/doi/10.1145/1721654.1721672>
- [2] W. Stallings, "Cryptography and Network Security: Principles and Practice," Pearson, 2017. Available at: <https://www.cs.vsb.cz/ochodkova/courses/kpb/cryptography-and-network-security-principles-and-practice-7th-global-edition.pdf>
- [3] R. Cattell, "Scalable SQL and NoSQL Data Stores," ACM SIGMOD Record, 2011. Available at: <https://dl.acm.org/doi/10.1145/1978915.1978919>
- [4] M. Bostock, V. Ogievetsky, and J. Heer, "D³ Data-Driven Documents," IEEE Transactions on Visualization and Computer Graphics, 2011. Available at: <https://ieeexplore.ieee.org/abstract/document/6064996>

- [5] S. Bird, E. Klein, and E. Loper, Natural Language Processing with Python, O'Reilly Media, 2009. Available at: <https://www.oreilly.com/library/view/natural-language-processing/9780596803346/>
- [6] F. Hutter, L. Kotthoff, and J. Vanschoren, Automated Machine Learning: Methods, Systems, Challenges, Springer, 2019. Available at: <https://link.springer.com/book/10.1007/978-3-030-05318-5>