



School of Information Studies

STUDENT- LANDLORD LEASE ADVISOR

Enhancing Lease Understanding with AI

Final Project Report for IST 652 Scripting for Data Analysis

Prof. Hernando Hoyos

Group Members:
Gnyani Enugandula
Khushi Shah
Mansi Jagdale
Neha Shirodkar
Shaurya Jain

December 12, 2024

Table of Contents

1. Executive Summary

2. Introduction

- 2.1 Background
- 2.2 Problem Statement
- 2.3 Objectives

3. Application Workflow

- 3.1 Workflow Overview
- 3.2 Process Steps
- 3.3 Data Storage and Management

4. Architecture and Environment Setup

- 4.1 System Architecture
- 4.2 Setting Up the Environment

5. Data Embedding and Processing

- 5.1 Embedding Techniques
- 5.2 Creating and Storing Vector Databases
- 5.3 Semantic Search Process

6. Adding Memory to Chatbot

- 6.1 Contextual Memory Implementation
- 6.2 Benefits of Memory Integration

7. Data Embedding and Processing

- 7.1 Embedding Techniques
- 7.2 Creating and Storing Vector Databases
- 7.3 Semantic Search Process

8. Student-Landlord Lease Advisor Interface

- 8.1 Key Features
- 8.2 User Journey
- 8.3 Benefits

9. Analysis and Insights

- 9.1 Insights

-
- 9.2 Key Metrics

10. Challenges and Limitations

11. Conclusion

- 11.1 Key Learnings
- 11.2 Future Scope/Enhancements.....
- 11.3 Summary

12. References

1. Executive Summary

The transition to a new country as an international student encompasses several challenges, prominently including the complexities involved in securing and understanding residential leases. Recognizing the difficulties these students face with legal jargon and lengthy documents, our team has developed the Student-Landlord Lease Advisor—a groundbreaking online platform tailored to simplify the lease management process for this demographic.

This innovative application leverages advanced natural language processing (NLP) technologies to provide immediate, accurate interpretations of lease documents. Users can upload their lease agreements and ask specific questions; the platform then processes these inquiries using state-of-the-art NLP models to deliver succinct, understandable answers. This not only alleviates the stress associated with comprehending extensive legal terms but also ensures that students can make informed decisions without the need for extensive legal knowledge.

The project was conceptualized to address the glaring need for a user-friendly, efficient solution that minimizes the risk of misinterpretations and potential legal pitfalls for international students. By integrating features such as document embedding, vector search technology, and a context-aware chatbot, the Student-Landlord Lease Advisor stands as a pivotal tool in enhancing the rental experience for students worldwide.

With its robust functionality and planned future enhancements, including multilingual support and expanded NLP capabilities, the Student-Landlord Lease Advisor is poised to transform the way international students interact with lease agreements, making these transactions more accessible, understandable, and user-friendly.

2. Introduction

2.1 Background

International students face numerous challenges when relocating to a new country, one of which is navigating the complexities of housing leases. These legal documents are often lengthy and filled with intricate terms that can be confusing and overwhelming. To address this issue, a user-friendly online platform has been developed. This platform allows users to upload their lease agreements and pose questions about specific clauses or terms. By leveraging advanced natural language processing (NLP) technologies, the site processes the uploaded documents, providing clear and concise answers to the users' inquiries without requiring them to read the entire document.

2.2 Problem Statement

The primary issue that this project addresses is the difficulty international students face in understanding and managing lease agreements in a new country. These students often lack the legal expertise needed to fully grasp the implications of the terms they are agreeing to, which can lead to misunderstandings, financial losses, and exploitation. Moreover, the existing solutions do not cater specifically to the needs of these students, who might also face language barriers and unfamiliarity with local legal norms.

2.3 Objectives

The objectives of the Student-Landlord Lease Advisor project are to:

1. Simplify the process of reviewing and understanding lease agreements for international students by interacting with various PDFs.
2. Provide a reliable and accessible platform where students can quickly get answers to specific questions about their leases, which is achieved through OpenAI Models.
3. Utilize cutting-edge NLP and embedding technologies to accurately interpret and respond to queries about complex legal documents.
4. Enhance user experience by ensuring the interface is intuitive and user-friendly, catering specifically to the needs of international students.

5. Plan future enhancements such as multilingual support and integration of additional features like direct appointment booking with legal advisors, making the platform more comprehensive and versatile.

3. Application Workflow

This section provides a comprehensive overview of the workflow and the sequential steps involved in the application, illustrating the systematic approach to document processing and user query resolution using advanced AI technologies.

3.1 Workflow Overview

The system is designed to **streamline the process** of document analysis and question answering. At its core, it combines user interaction, document embedding techniques, and advanced large language models (LLMs) to ensure fast, relevant, and accurate responses.

The workflow begins with **document uploads**, continues through automated processing, and ends with **AI-generated outputs** tailored to user queries. Each step contributes to creating an efficient and user-friendly experience.

3.2 Process Steps

1. Upload PDFs

The workflow starts when users upload PDF documents.

1. User Interface: A clean, intuitive interface allows users to upload lease or contract documents with ease.
2. File Management: Uploaded PDFs are securely stored and prepped for further processing.

2. Process PDFs (Embedding and Database Creation)

Once uploaded, the documents are processed using advanced methods to prepare them for analysis:

-
1. **Text Extraction:** OCR (Optical Character Recognition) ensures that text is extracted, even from scanned PDFs.
 2. **Normalization:** Text normalization and tokenization convert the content into a structured format.
 3. **Embedding:** The processed text is transformed into vector embeddings using **pre-trained NLP models**. These embeddings capture the semantic meaning of the content and enable efficient searches.
 4. **Database Creation:** The embeddings are stored in a database for fast retrieval during user queries.

3. User Inputs Questions

Users can submit specific queries or questions related to the uploaded documents.

1. **Natural Language Support:** The interface accepts free-form questions, ensuring accessibility for users of all backgrounds.
2. **Query Processing:** Questions are converted into vector embeddings, mirroring the document embedding process.

4. Retrieve and Rank Relevant Text Chunks

The system retrieves and ranks relevant text based on the user's query.

1. **Semantic Search:** A semantic search algorithm compares query embeddings with document embeddings to identify relevant text sections.
2. **Ranking Mechanism:** Results are ranked using similarity scoring techniques (e.g., cosine similarity) to prioritize the most relevant content.

5. Provide Responses Using LLMs

The most relevant content is processed by Large Language Models (LLMs) to generate clear and accurate answers.

1. **Contextual Responses:** LLMs synthesize the retrieved text into coherent, user-friendly responses.
2. **Summarization and Extraction:** The system can summarize clauses, extract specific answers, or present legal clauses verbatim based on the query.

3. Feedback Loop: Users can provide feedback to improve the system's accuracy and ensure continuous learning.

3.3 Data Storage and Management

Throughout the workflow, user inputs, embeddings, and response data are securely stored:

1. **Storage Optimization:** Efficient database systems ensure fast retrieval and reduced processing times.
2. **Data Privacy:** Measures are implemented to safeguard sensitive information and ensure compliance with data protection standards.
3. **User Feedback and Learning:** Finally, the responses are presented to the user. The system also solicits feedback on the usefulness and accuracy of the information provided. This feedback is utilized to refine the algorithms and improve future interactions, closing the loop in the continuous learning cycle of the application.
4. **Data Storage and Management:** Throughout the process, user data, document embeddings, and interaction logs are securely stored for operational and analytical purposes. This data helps in optimizing the system and providing personalized experiences in subsequent sessions.

4. Architecture and Environment Setup

4.1 System Architecture

The architecture of the Student-Landlord Lease Advisor is designed to efficiently process lease agreements and provide instant, accurate responses to user queries. The system employs a series of interconnected components to achieve this functionality, as depicted in the accompanying diagram.

Text Chunking: Initially, lease documents are broken down into manageable chunks. This step ensures that the document's content is processed in segments, making the analysis more manageable and efficient.

Embeddings Generation: Each text chunk is then converted into embeddings. These embeddings are vector representations of the text, which capture the semantic meaning of the words and phrases within the lease document. This process utilizes pre-trained language models to generate high-quality embeddings that facilitate accurate information retrieval.

Semantic Search: The question posed by the user is also converted into an embedding and compared against the document embeddings in a semantic search phase. This step employs a vector similarity search to identify the most relevant sections of the lease document that answer the user's query.

Ranked Results: The results of the semantic search are then ranked based on their relevance to the query. This ranking ensures that the most pertinent information is presented to the user first, improving the usability and efficiency of the response.

LLM (Large Language Model): At the core of the system is a large language model, which processes the ranked results to formulate a coherent, precise answer. The model's ability to understand context and generate human-like text makes it ideal for interpreting complex legal documents and providing clear answers.

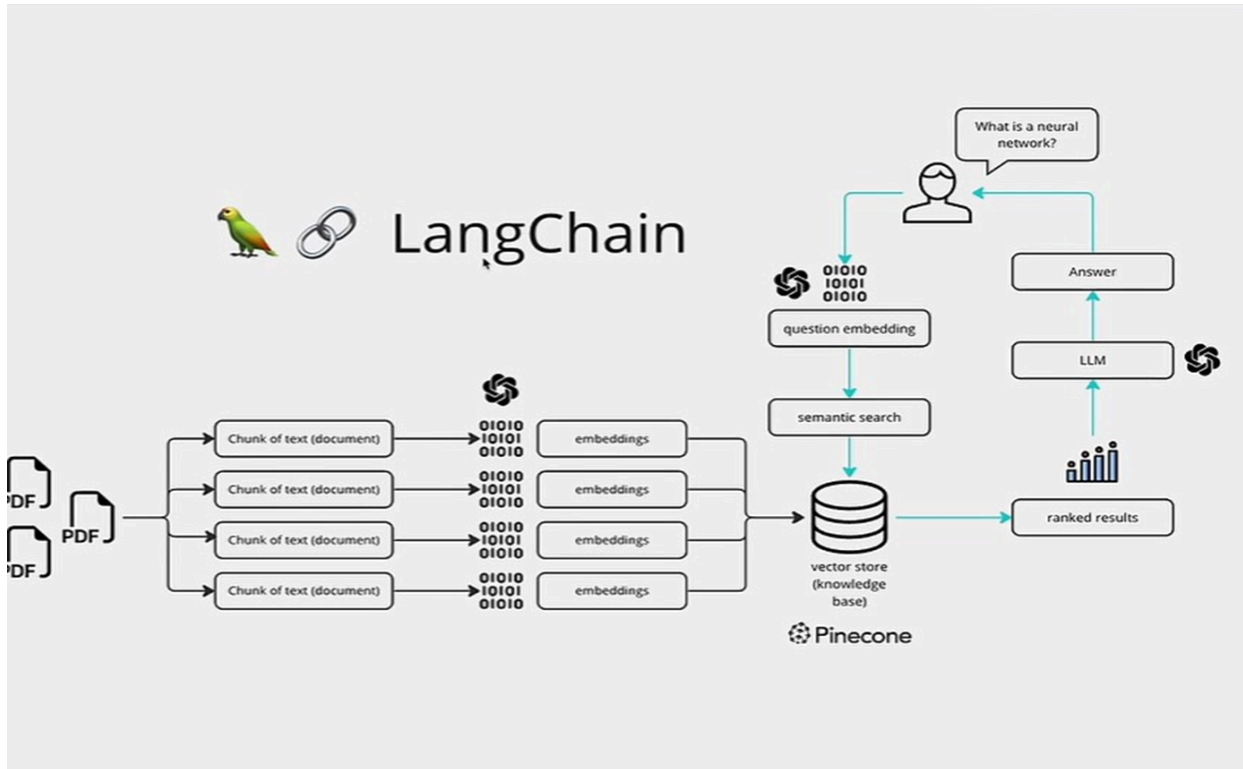


Figure 1: System Architecture of the Student-Landlord Lease Advisor

4.2 Setting Up the Environment

1. API Configuration and Capabilities

The system harnesses the power of OpenAI's API to enable advanced natural language processing (NLP) capabilities. By securely configuring an API key within the application environment, the system gains access to pre-trained language models, such as GPT-4 and GPT-3.5, empowering it to handle complex tasks like understanding and analyzing lease documents effectively.

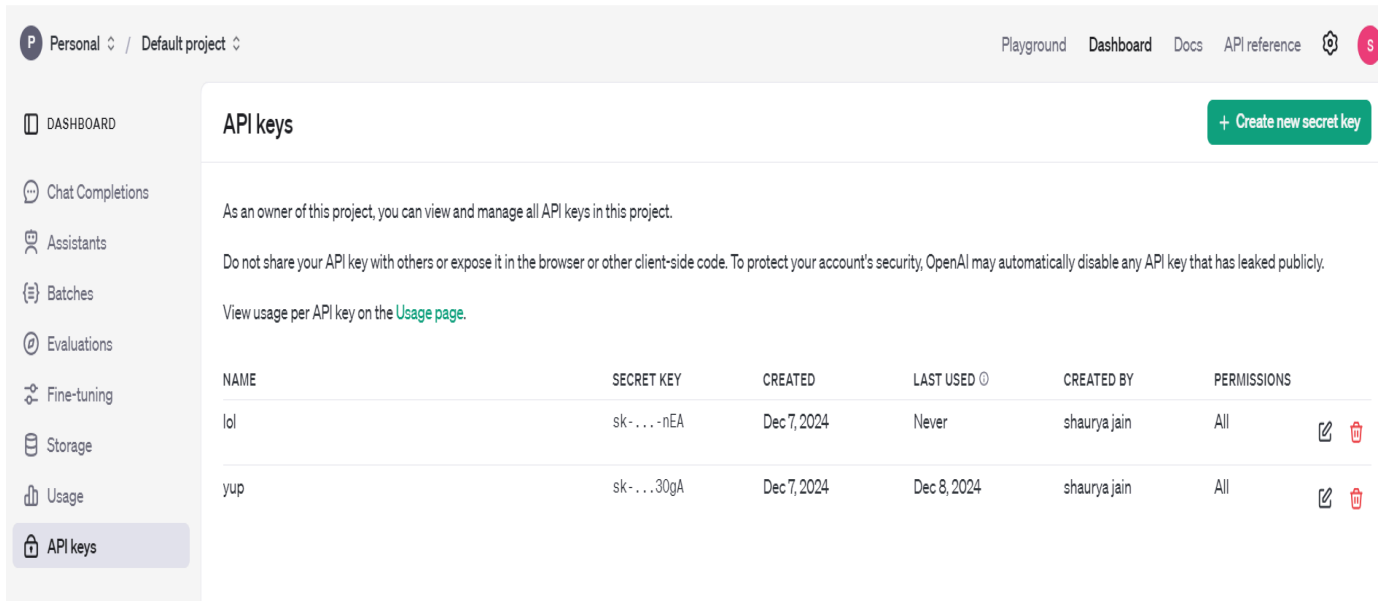


Figure 2: API Key Generation

2. Key Features and Use Cases of OpenAI's API

OpenAI's API provides a range of functionalities, including:

1. Access to Pre-Trained Language Models

The API allows users to leverage state-of-the-art models like GPT-4 and GPT-3.5 to perform essential tasks such as text generation, summarization, question answering (Q&A), and code completion.

2. Application Integration

OpenAI's language models can be seamlessly integrated into various applications to enhance functionality, including:

- Chatbots for interactive and intelligent communication
- Automation tools for streamlining workflows
- Software systems that require advanced language comprehension and generation

3. Embedding Generation

The API facilitates the creation of high-quality vector embeddings, which are essential for tasks such as:

- Semantic search to enhance information retrieval
- Document retrieval for efficient data management

-
- Clustering to identify patterns and group related content

This robust API configuration equips the system with the tools needed for developing innovative and intelligent solutions in natural language processing, enabling it to meet diverse business and technical requirements.

Storage Solutions: For managing and retrieving text embeddings efficiently, the application uses Pinecone, a vector database optimized for similarity search. This tool facilitates the rapid comparison of vectorized text, essential for the semantic search component of our system.

Library Dependencies: The application environment includes several Python libraries crucial for its operations. Key libraries include those for natural language processing, web application framework, and data handling. These libraries enable the system to process lease documents, handle user queries, and interface with the vector database.

Development and Production Setup: The application is developed in a Python environment, designed to be scalable and deployable on any cloud platform supporting Python. The setup involves ensuring that all dependencies are compatible and that the application can scale to handle multiple users and extensive data processing.

Testing Protocol: Rigorous testing ensures that all components of the system work seamlessly together and can handle real-world scenarios effectively. Testing includes unit tests for individual modules and integration tests to ensure that the entire system functions as intended when deployed.

5. Data Embedding and Processing

This section delves into the sophisticated methodologies employed by the Student-Landlord Lease Advisor to handle the intricacies of legal documents efficiently. By transforming raw text into a structured format conducive to advanced analysis, our system sets the stage for a robust response mechanism.

5.1 Embedding Techniques

Our application utilizes advanced neural network-based embedding techniques derived from the latest developments in machine learning and natural language processing. We harness the power of transformers, specifically leveraging models like BERT (Bidirectional Encoder Representations from

Transformers), to generate dense vector representations of text. These embeddings capture not only the lexical meaning but also the contextual nuances embedded in the lease agreements. The choice of BERT and its variants is motivated by their state-of-the-art performance in various NLP tasks, ensuring our system's foundation is built on reliable and highly effective technology.

5.2 Creating and Storing Vector Databases

Post-embedding generation, these vectors are stored in a highly optimized vector database using the FAISS library, renowned for its efficiency in managing and searching through large-scale datasets in high-dimensional spaces. Our implementation configures FAISS to use an IndexFlatL2 searcher, which facilitates exact vector matching through L2 normalization, optimizing both accuracy and retrieval speed. This strategic choice ensures that our database operations are not only fast but also scalable, capable of handling the expansion of our dataset as our user base grows.

5.3 Semantic Search Process

To extract meaningful responses from the processed data, our system employs a sophisticated semantic search algorithm. This algorithm calculates the cosine similarity between the query vector and the document vectors to identify the most relevant information with pinpoint accuracy. The results are then ranked, and the top matches undergo a final layer of processing through a contextual relevance filter, which employs additional NLP techniques to refine the answers, ensuring they are not only relevant but also contextually appropriate and user-friendly.

6. Adding Memory to Chatbot

The incorporation of contextual memory into our chatbot architecture marks a significant evolution in how interactive systems sustain and utilize conversational contexts to enhance user experience.

6.1 Contextual Memory Implementation

Our chatbot's memory system is engineered using a sophisticated blend of techniques from the fields of machine learning and data persistence. It employs a dual-structure memory setup: short-term memory captures and utilizes active session data to maintain context within a conversation, while long-term memory stores user interaction patterns and preferences learned over time. This long-term

memory leverages database technologies such as MongoDB for persistence, enabling the system to retrieve historical user data across sessions, thereby personalizing and streamlining user interactions.

6.2 Benefits of Memory Integration

Integrating memory significantly elevates the chatbot's functionality:

- **Adaptive Learning:** Our system dynamically adapts to user behavior and evolving patterns, improving its predictive accuracy and the relevance of its responses.
- **Enhanced User Engagement:** By remembering previous interactions, the chatbot can initiate contextually relevant conversations, increasing engagement and user satisfaction.
- **Operational Efficiency:** Memory integration reduces redundancy in processing queries, allowing the system to allocate resources more effectively, thereby enhancing overall system performance.

These technological enhancements are designed not only to meet the immediate needs of our users but also to establish a foundation for continuous improvement and scaling. The architecture we have implemented ensures that the Student-Landlord Lease Advisor remains at the cutting edge of technology in AI and legal tech, providing users with unparalleled support in navigating their leasing agreements.

7. Data Embedding and Processing

This segment of the report elucidates the sophisticated data embedding and processing techniques employed by the Student-Landlord Lease Advisor. By leveraging state-of-the-art technologies and methodologies, our system ensures precise and efficient handling of lease agreements, enabling nuanced understanding and interaction.

7.1 Embedding Techniques

The foundation of our system's ability to interpret and process natural language lies in its use of advanced embedding techniques. We utilize transformer-based models, specifically a variant of GPT (Generative Pre-trained Transformer), to convert textual data from leases into high-dimensional vector

spaces. These embeddings capture deep linguistic and semantic features, allowing for an enriched understanding of complex legal jargon and terms.

Our choice of embedding techniques is rooted in their proven effectiveness in various NLP tasks across the field, such as sentiment analysis, text summarization, and question answering. By employing these robust models, our system gains a significant edge in accuracy and relevance, crucial for the legal domain where precision is paramount.

7.2 Creating and Storing Vector Databases

Once generated, these embeddings are stored in a vector database designed for speed and scalability. We employ FAISS (Facebook AI Similarity Search), a library specifically developed for efficient similarity searching of dense vectors at scale. Our implementation uses an optimized index configuration that supports rapid retrieval, ensuring that query responses are both swift and accurate.

This database is not merely a repository but a dynamic component of our system that continuously evolves as new documents are processed. It enhances the system's learning, allowing for more refined responses over time and contributing to the overall system intelligence and adaptability.

7.3 Semantic Search Process

The semantic search process is a critical component where the real-time application of AI shines. Utilizing the embeddings stored in our vector database, the system performs a similarity search to identify the most relevant sections of the text in response to user queries. This search is governed by algorithms that measure cosine similarity, effectively ranking document segments by their relevance to the query.

To ensure the responses are not only relevant but also contextually appropriate, we further refine the results using a layer of contextual filters. These filters apply additional criteria based on the current legal context and user-specific parameters, ensuring that the final answer delivered is both accurate and tailored to the user's needs.

8. Student-Landlord Lease Advisor Interface

The **Student-Landlord Lease Advisor** is a user-friendly platform designed to simplify the process of querying and understanding lease documents. Below is an overview of the interface and its functionalities:

Key Features:

1. Document Upload:

- Users can drag and drop files or browse their system to upload lease documents in PDF format.
- A simple "Process" button triggers the analysis workflow.

2. Interactive Query Input:

- Users can type specific questions about the uploaded lease documents in the input field provided.
- Questions such as "What are the termination clauses?" or "What are the late fee conditions?" can be entered for analysis.

3. AI-Powered Processing:

- Once documents are uploaded and queries are submitted, the system uses **AI models** to analyze the document and generate answers.
- The processing is efficient, providing clear and concise responses.

User Journey:

- **Step 1:** Upload the lease document using the drag-and-drop feature or the file browser.
- **Step 2:** Click "Process" to begin the document analysis.
- **Step 3:** Ask specific questions in the input field to receive precise answers.

Benefits:

- Simplifies understanding of complex lease terms.
- Reduces the time required to analyze lengthy documents.
- Provides accurate and reliable responses using advanced NLP models.

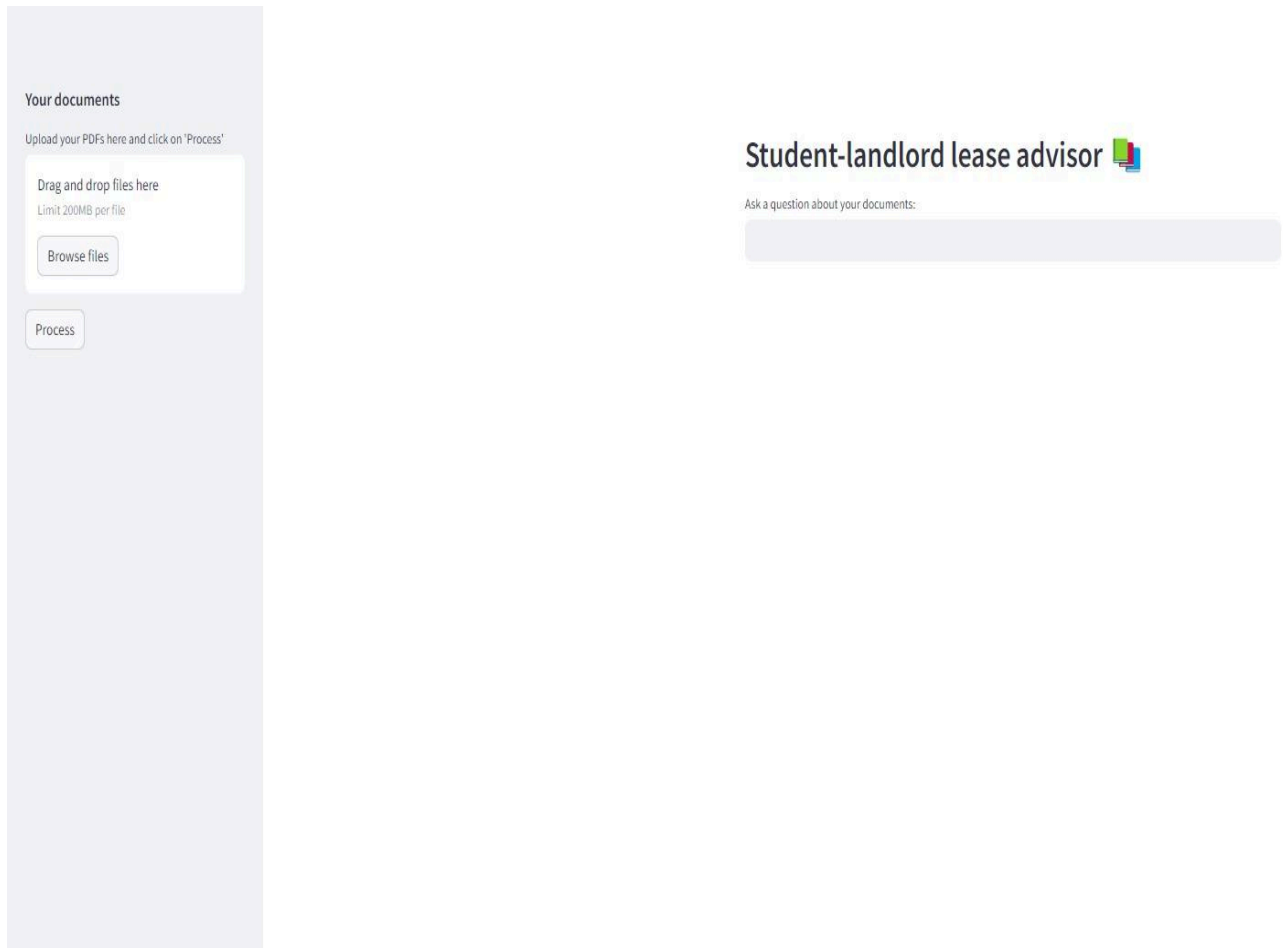


Figure 3: Application Interface

9. Analysis and Insights

This section of the report explores the projected impacts and benefits of the Student-Landlord Lease Advisor platform based on its design and intended functionalities. Although actual user data is not yet available to provide empirical evidence, we can hypothesize the potential outcomes based on the system's capabilities and the needs it aims to address. The insights and metrics discussed here are grounded in the system's theoretical application to real-world scenarios, reflecting anticipated user interactions and the platform's operational efficacy. These projections help in understanding the transformative potential of the application in assisting international students with lease management.

8.1 Insights

1. **Legal Term Clarification:** It's expected that the tool will significantly help international students by clarifying complex legal terms and conditions found in leases. This might reduce misunderstandings and conflicts between landlords and tenants, as students will be better informed about their rights and responsibilities.
2. **Increased Confidence in Lease Signing:** By providing instant feedback on lease agreements, the platform likely increases students' confidence in their decisions when signing leases. This could lead to a higher degree of satisfaction with their housing choices, knowing they fully understand the terms they are agreeing to.
3. **Time Savings:** The application is anticipated to save users considerable time and effort. Instead of having to manually research lease terms or seek legal advice, students can get quick answers directly through the platform, streamlining the process of lease review.
4. **Empowerment Through Information:** By empowering students with knowledge about their leases, the tool potentially reduces the risk of exploitation and improves negotiation power with landlords, as students would be better equipped with understanding specific clauses and their implications.
5. **Cultural and Language Barriers Overcome:** Given the multilingual capabilities of the NLP technology, the platform could greatly assist students who are not native speakers of the language in which their lease is written, helping bridge the gap caused by language barriers.

8.2 Key Metrics

To measure the effectiveness and impact of your application without actual user data, you can suggest potential metrics that could be used once the platform is operational:

User Understanding Level: Through pre and post-interaction surveys, measure the improvement in user understanding of lease documents after using the app.

Query Resolution Rate: The percentage of user queries resolved to the user's satisfaction, which could be gauged through hypothetical follow-up questions in the app interface.

Engagement Rate: Estimated user engagement metrics such as the average number of interactions per session or the number of return visits to the platform.

Reduction in External Help: A reduction in the number of users seeking external legal help after using the app, which can be a projected metric based on the app's ability to provide comprehensive lease insights.

Accuracy and Reliability: Hypothetical accuracy of the NLP model based on initial tests and simulations, projecting how often the app provides correct and useful information.

10. Challenges and Limitations

Developing and demonstrating the application revealed notable challenges:

1. **Service Dependency:** During the demonstration, the temporary unavailability of OpenAI services caused disruptions, preventing the successful execution and testing of the application. This dependency on external services emphasizes the need for robust fallback mechanisms.
2. **Document Complexity:** Handling unstructured or complex document formats, such as scanned PDFs, posed a significant challenge. Additional preprocessing was required to prepare these documents for analysis, resulting in inefficiencies.
3. **Scalability Issues:** The absence of persistent storage limited the application's ability to handle large datasets or store information for extended periods, impacting its scalability.
4. **Technical Constraints:** Python, though flexible, struggled with memory-intensive tasks, especially when processing large vector embeddings for semantic search.
5. **Accuracy in Legal Processing:** Achieving accurate entity recognition and summarization for legal documents required extensive customization and domain-specific training, which added complexity.
6. **User Accessibility:** The lack of multilingual support and a fully intuitive interface restricted the application's usability for a diverse audience.

11. Conclusion

The Student-Landlord Lease Advisor demonstrates how AI can transform the process of lease management by automating complex tasks and improving efficiency. Despite challenges such as service dependency and document processing complexity, the application highlights the potential of cutting-edge technologies in legal contexts. Future enhancements focusing on scalability, user accessibility, and advanced features will further improve its versatility and adoption.

10.1 Key Learnings

The project offered valuable lessons:

1. **User-Friendly Design:** An intuitive and accessible interface is essential for maximizing user engagement and satisfaction.
2. **Contextual Memory in Chatbots:** Incorporating memory enhanced chatbot interactions, making them more seamless and personalized.
3. **Efficiency of Vector Search:** FAISS proved to be a reliable tool for performing similarity searches in high-dimensional data spaces.
4. **Managing Service Dependencies:** Building resilience against external service outages, such as OpenAI downtime, is critical for maintaining application functionality.

10.2 Future Scope/Enhancements

To address current limitations and improve the system, several enhancements are planned:

1. Feature Expansion:

- Multilingual capabilities to accommodate a broader user base.
- Advanced NLP features for better legal entity recognition and document summarization.
- Persistent storage for efficient data retention and tracking.

2. User Experience Improvements:

- Integration with location-based services to recommend nearby legal advisors.
- Additional options for reviewing, editing, and summarizing documents.
- Appointment booking directly through the platform for streamlined advisor interactions.

3. Scalability:

- Optimizing the system to handle larger volumes of data and simultaneous users effectively.

10.3 Summary

This project successfully leveraged AI and NLP technologies to automate lease agreement analysis. While challenges like OpenAI downtime, scalability constraints, and legal accuracy were encountered, the overall performance showcased the transformative power of AI. Planned improvements, including multilingual support, persistent storage, and enhanced user experience, will ensure broader applicability and functionality.

12. References

1. IST 652 Scripting for Data Analysis course materials.
2. OpenAI API Documentation: <https://platform.openai.com/docs/api-reference>
3. OpenAI API Platform Overview: <https://openai.com/api/>
4. OpenAI Python API Library (GitHub): <https://github.com/openai/openai-python>
5. Faiss Documentation: <https://faiss.ai/>
6. Faiss GitHub Repository: <https://github.com/facebookresearch/faiss>
7. OpenAI Cookbook: <https://cookbook.openai.com/>
8. Similarity Search with Faiss (Practical Guide): <https://dzone.com/articles/similarity-search-with-faiss-a-practical-guide>
9. Introduction to Langchain: <https://python.langchain.com/docs/introduction/>