

Name Qin Du, Ziyu Lin, Lu Huang Date: 2024.12.19

(last name, first name)

ID: N10961564, N17172561, N11560862

Course Section Number: CSCI-GA.2433-001

Project Part 4

Total in points (100 points total): _____

Professor's Comments:

Affirmation of my Independent Effort Qin Du, Ziyu Lin, Lu Huang
(Sign here)

Github link: https://github.com/ThisZiyu/DBMS_2024FALL.git

1. Project Background

In today's data-driven environment, enterprises face significant challenges in integrating structured and unstructured data to derive real-time insights. Such integration is crucial for achieving operational excellence and maintaining a competitive edge. This project aims to design and implement an end-to-end insurance management system that integrates machine learning (ML) models, databases, and real-time data pipelines. The proposed solution addresses the need for personalized insurance pricing, automated claims processing, and dynamic risk assessment while prioritizing fairness and transparency.

2. DIKW Reference Architecture

The DIKW (Data, Information, Knowledge, Wisdom) framework provides the foundational architecture for the project. This hierarchy transforms raw data into actionable insights. Health-related data, such as demographic, lifestyle, and medical factors collected from `standardized_health_data.csv`, serves as the foundational input. These structured datasets guide risk predictions and insurance pricing decisions, with insights from ML models categorizing risks and informing premium calculations. The implementation of real-time premium updates and claim status tracking further enhances user experience and organizational efficiency. Governance strategies include ETL processes to ensure data quality, logical schemas for metadata management, and Random Forest model evaluations to identify and mitigate potential biases.

3. End-to-End Solution Design

This project is designed around three key business use cases.

1. The first use case focuses on insurance quotation, where customers input demographic and health data such as age, gender, and BMI. A Random Forest model, `health_risk_model.pkl`, evaluates these inputs to predict health risk levels categorized as Low, Medium, or High. Premiums for various insurance plans are dynamically calculated using base premiums adjusted for age groups and risk multipliers specific to each risk category.
2. The second use case involves claims management. Policyholders submit claims that include claim numbers, policy details, and service descriptions. The system generates unique claim IDs, tracks statuses such as Submitted, Approved, or Rejected, and provides APIs for real-time updates.
3. The third use case pertains to real-time premium updates, where premiums are adjusted dynamically based on market conditions, and timestamped updates are displayed on user dashboards. Market adjustments refer to dynamic modifications in insurance premiums based on real-time external factors. These adjustments are implemented to reflect:

- **Economic Indicators:** Fluctuations in inflation rates or healthcare costs may lead to percentage-based increases or decreases in premium rates.
- **Policy Demand and Supply Trends:** Increased demand for certain types of insurance coverage may trigger slight upward adjustments, whereas reduced demand could result in discounts.
- **Regulatory Changes:** Compliance with updated government policies or industry regulations may necessitate premium recalibrations.
- **Competitor Benchmarking:** Real-time monitoring of competitor pricing strategies ensures competitiveness by adjusting premiums accordingly.

The system incorporates these adjustments as percentage modifiers ($\pm 2\%$) applied to the base premium calculations. The dynamic updates are automated via backend scripts that pull relevant data feeds, evaluate the external conditions, and apply the adjustments in near real-time. Users see these adjustments reflected instantly on their dashboards, ensuring transparency and alignment with current market dynamics.

The system architecture integrates a web-based frontend for user interactions, a Flask-based backend for data processing and ML model invocation, an SQLite database managed via SQLAlchemy ORM, and a Random Forest classifier for health risk predictions. The workflow begins with user data inputs via the web interface, followed by processing through the Flask backend. The ML model evaluates health risks, stores calculated premiums and risk levels in the database, and displays results to users. Claims are submitted and processed with real-time tracking of their statuses.



Figure 1: Insurance System Class Diagram (Business Use Cases)

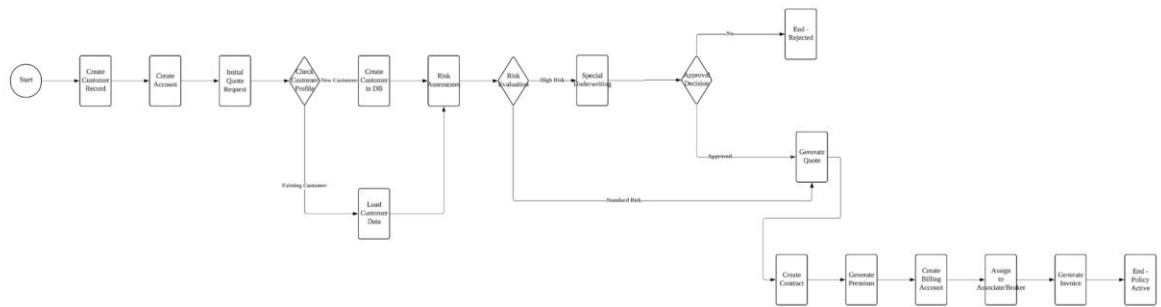


Figure 2: Business Workflow Diagram

The following diagram illustrates the end-to-end business workflow for the insurance management system, showcasing the integration of various processes from customer onboarding to policy activation. This workflow aligns with the three key business use cases described in the solution design:

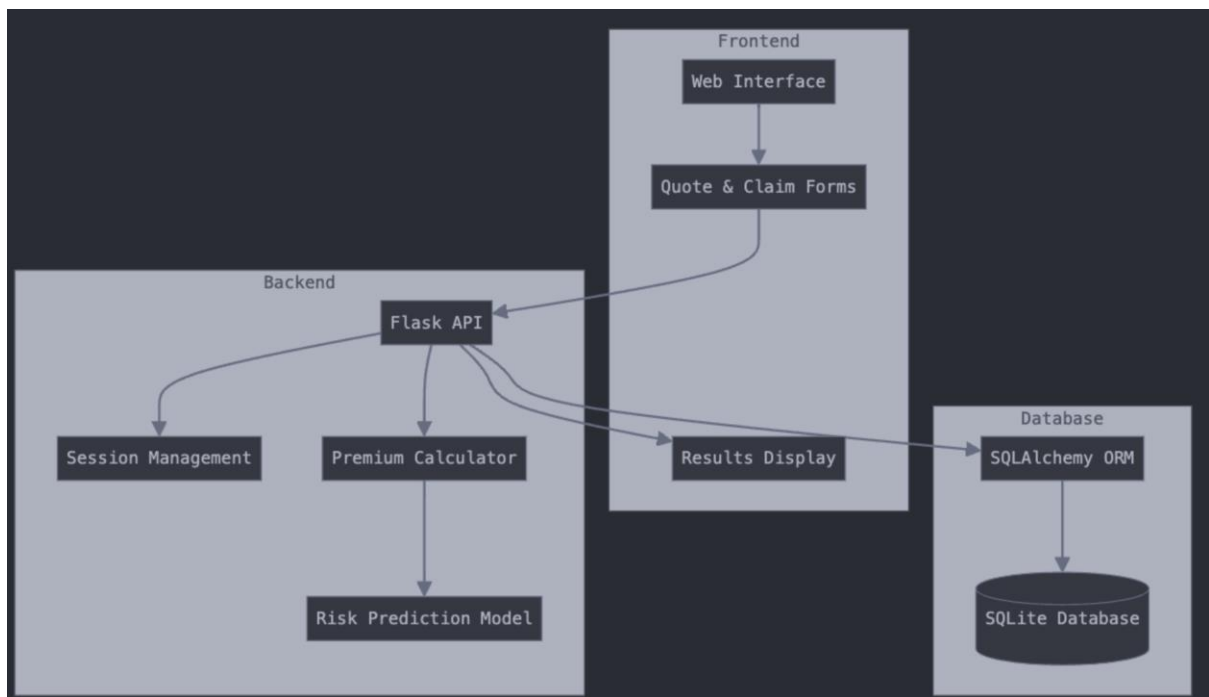
- **Customer Onboarding and Account Creation:** The process begins with customer profiling, where new customer records are created, and existing customer data is retrieved for further processing.
- **Risk Assessment and Evaluation:** Customers' health and demographic data are evaluated using a machine learning model. Based on the risk level (Low, Medium, High), further actions such as special underwriting or approval decisions are triggered.
- **Insurance Quotation and Policy Issuance:** Approved applications lead to the generation of insurance quotes, followed by the creation of contracts, premium calculations, and billing setup.
- **Policy Activation:** Once billing and contract processes are complete, the policy is activated and associated with a broker or agent for customer support.
- This workflow visually represents the system's capability to manage the entire insurance lifecycle, ensuring a seamless and automated experience for users.

4. Architecture Diagram

The architecture diagram of the system illustrates the integration of its core components:

- **Frontend:** A web interface where users input their data and view results.
- **Backend:** A Flask API that handles data processing, invokes the ML model, and manages user sessions.
- **Database:** An SQLite database, optimized through SQLAlchemy ORM, storing structured information for rapid retrieval.
- **Machine Learning Module:** A Random Forest model deployed to predict health risks and dynamically calculate insurance premiums.

The workflow depicted in the diagram includes data flow from user input to real-time processing and output visualization.



5. Data-Driven Module

5.1 Dataset Overview

- **Source:** standardized_health_data.csv
- **Features:** Includes age group, gender, BMI, smoking habits, alcohol consumption, education level, income, and pre-existing conditions.
- **Volume:** 5,000 synthetic records representing chronic disease trends.

	Age_Group	Gender	Race	Smoking	Alcohol_Consumption	\
0	31-45	Male	Caucasian	Never	Occasional	
1	61+	Male	Caucasian	Former	Never	
2	46-60	Female	Caucasian	Never	Never	
3	46-60	Male	Hispanic	Never	Occasional	
4	18-30	Female	Caucasian	Never	Occasional	

	Education_Level	Insurance_Type	Exercise_Frequency	Diet_Quality	\
0	Less than high school	Uninsured	1-2 times/week	Fair	
1	Less than high school	Private	Never	Poor	
2	Some college	Medicaid	1-2 times/week	Poor	
3	Graduate degree	Private	Never	Poor	
4	Some college	Uninsured	3-4 times/week	Fair	

	Sleep_Hours	Stress_Level	Occupation	Family_History	BMI_Category	\
0	6-7	High	Manual Labor	No	Normal	
1	7-8	Medium	Manual Labor	No	Normal	
2	8+	Low	Service Industry	No	Overweight	
3	7-8	Medium	Office/Desk Job	Yes	Normal	
4	6-7	Medium	Service Industry	Yes	Overweight	

	Income_Level	Disease_Type
0	Less than \$30,000	Colorectal Cancer
1	\$75,000-\$100,000	Diabetes
2	\$30,000-\$50,000	Diabetes
3	\$30,000-\$50,000	Diabetes
4	\$30,000-\$50,000	Sleep Apnea
		Chronic Respiratory Disease
		Sleep Apnea
		Obesity
		Lung Cancer
		Diabetes
		Colorectal Cancer
		Hypertension
		Heart Disease
		Depression
		Arthritis

Name: Disease_Type, dtype: int64

5.2 Machine Learning Model

A Random Forest classifier trained on this dataset predicts health risks with an overall accuracy of 62.4%. The model's performance metrics reveal high precision and recall for Low and Medium Risk categories, though High Risk predictions exhibit moderate recall and precision. Feature importance analysis highlights the influence of diet quality, stress levels, and age group on risk predictions.

The trained model, deployed via Flask, dynamically calculates premiums, while ETL pipelines monitor unstructured data changes and trigger model re-training when necessary thresholds are exceeded.

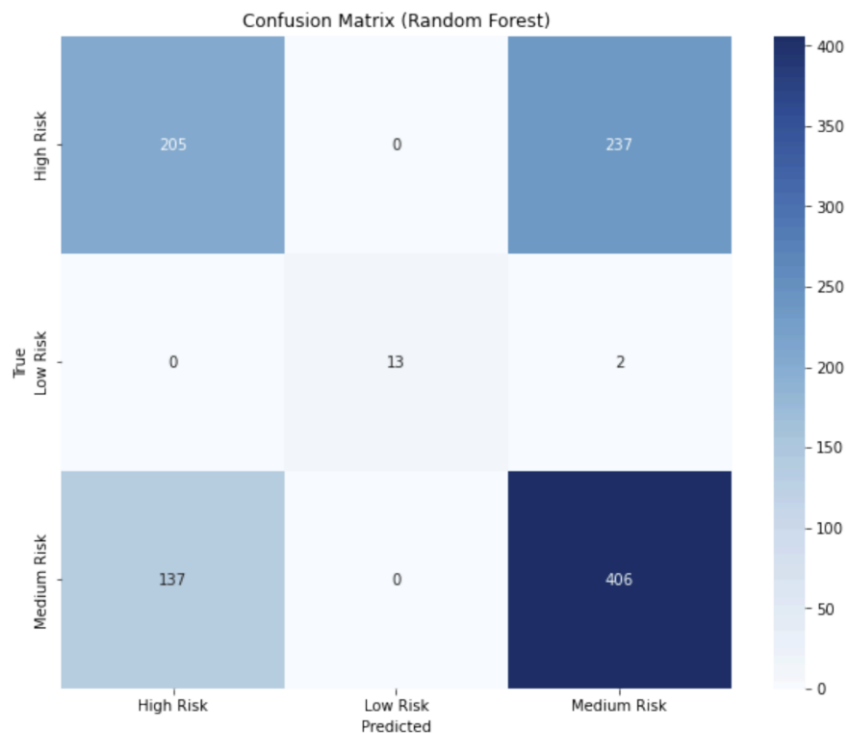


Figure 3: Confusion Matrix

Random Forest Test Accuracy: 0.624

Classification Report:

	precision	recall	f1-score	support
High Risk	0.60	0.46	0.52	442
Low Risk	1.00	0.87	0.93	15
Medium Risk	0.63	0.75	0.68	543
accuracy			0.62	1000
macro avg	0.74	0.69	0.71	1000
weighted avg	0.62	0.62	0.62	1000

Figure 4: Model performance

5.3 Model Deployment

The trained model (health_risk_model.pkl) was deployed using Flask, enabling real-time predictions. The system dynamically integrates predictions into insurance premium calculations, applying risk multipliers (0.8x for Low Risk, 1.0x for Medium Risk, and 1.5x for High Risk).

5.4 Automation and Re-Training

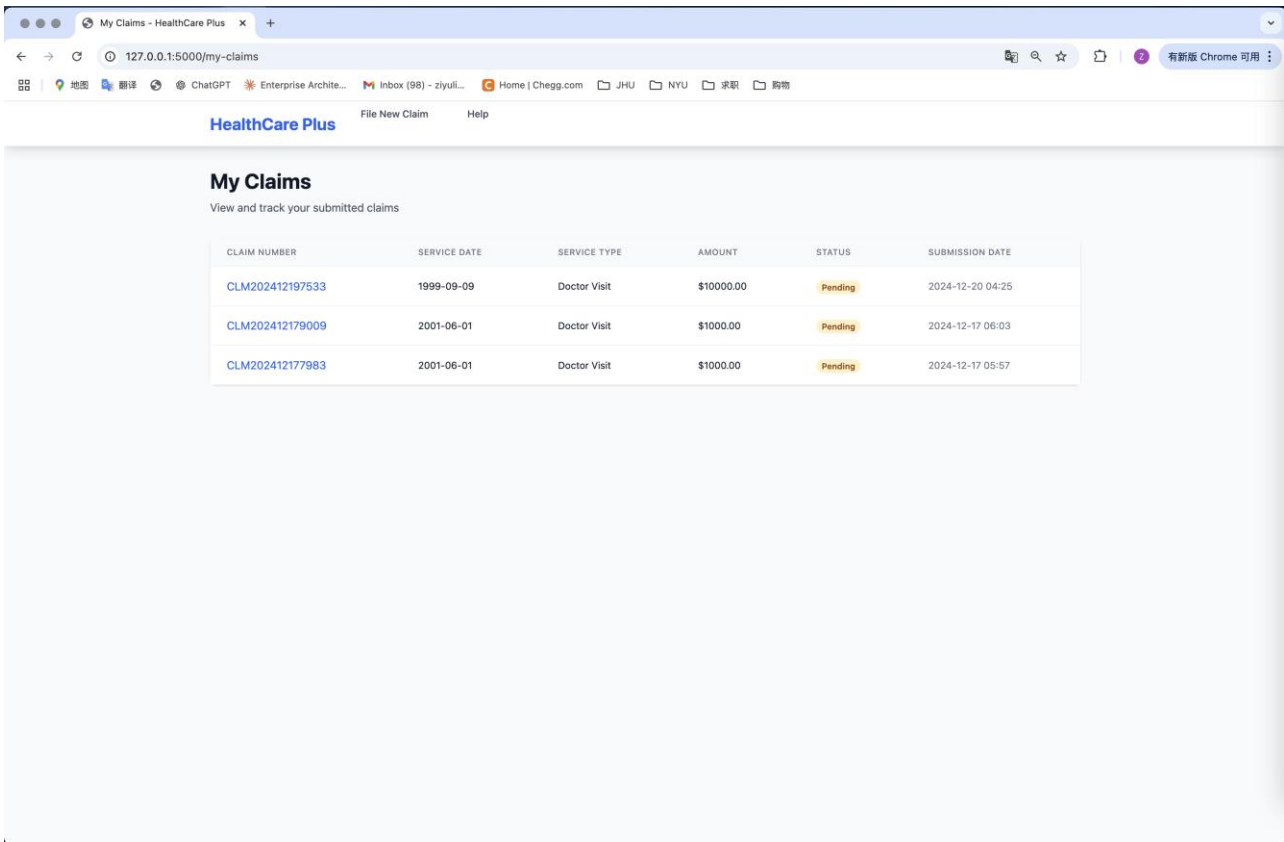
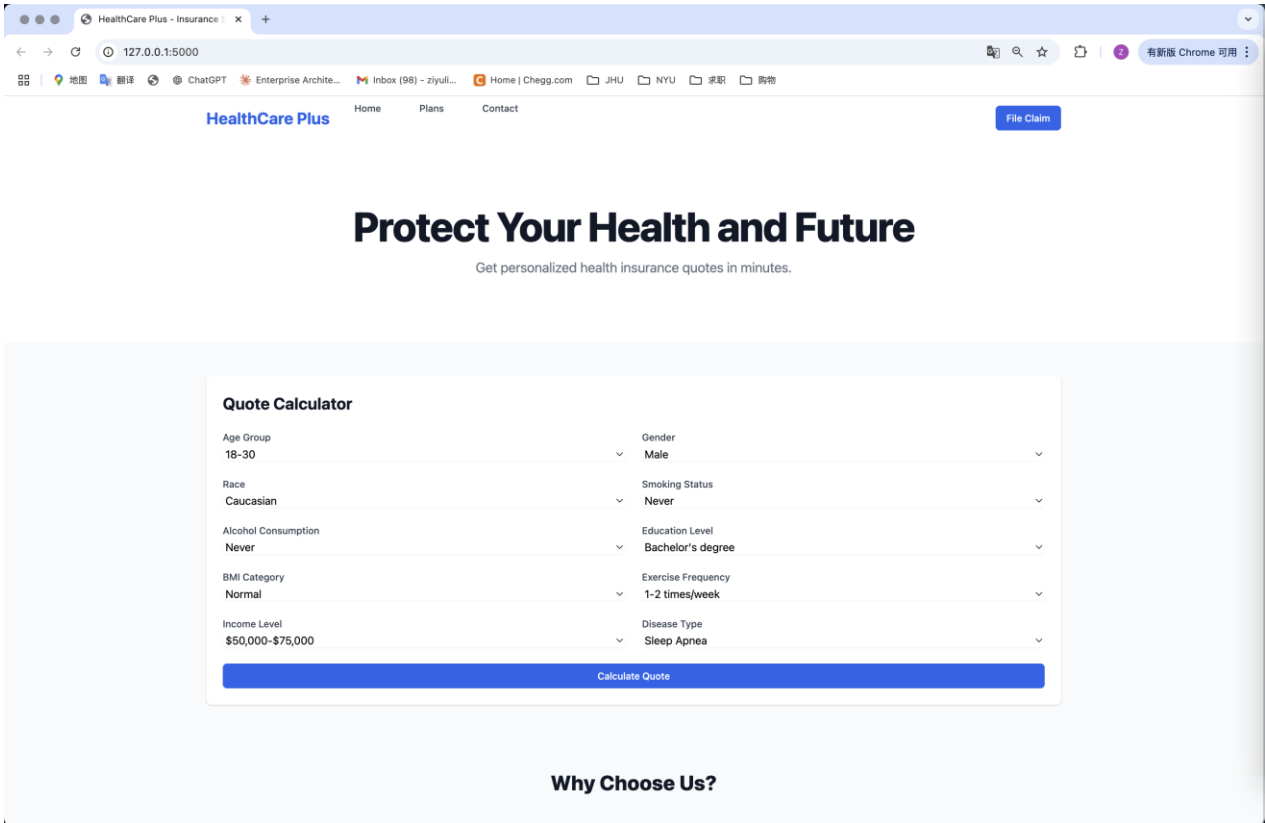
- **Pipeline:** ETL monitors unstructured data changes.
- **Triggers:** Data shifts exceeding thresholds initiate model re-training.
- **Tools:** Joblib for model serialization; Pandas for preprocessing.

6. Backend and API Functionality

The backend of the system, developed using Flask, integrates seamlessly with an SQLite database and performs critical operations that drive the system's functionalities. Risk prediction is managed through a Random Forest model (health_risk_model.pkl), which evaluates user input to categorize risk levels into Low, Medium, or High. The predict_risk function formats user data for the model and handles cases where the model is unavailable by defaulting to Medium Risk predictions. Premium calculation follows a structured approach with predefined base rates for different age groups: \$100 for ages 18-30, \$150 for 31-45, \$200 for 46-60, and \$250 for those 61 and older. These base premiums are adjusted based on risk levels, applying multipliers such as 0.8x for Low Risk, 1.0x for Medium Risk, and 1.5x for High Risk. Additional factors, like the selected insurance plan (Basic, Standard, Premium), further influence premiums by modifying coverage and deductibles.

The system facilitates user interactions through several RESTful API endpoints, ensuring real-time communication between the user interface, database, and ML model. The /calculate_quote endpoint accepts demographic and health data to calculate premiums, while /submit_claim allows policyholders to submit insurance claims for database storage. The /claim_status/<claim_id> endpoint retrieves the current status of submitted claims, providing timely updates to users. All user data, including claims and quotes, is stored in an SQLite database managed via SQLAlchemy ORM, ensuring efficient data retrieval and real-time updates for dynamic premium adjustments.

To maintain system reliability, robust error-handling mechanisms are implemented. For example, if the ML model fails to load, the system defaults to Medium Risk predictions to prevent interruptions. These functionalities, combined with seamless database integration and real-time processing capabilities, enable a comprehensive and user-friendly insurance management system.



File a Claim - HealthCare Plus

127.0.0.1:5000/file-claim

有新版 Chrome 可用

HealthCare Plus

My Claims

Help

File a New Claim

Submit your medical claim for reimbursement

Patient Information

Patient Name

Ziyu

Policy Number

123456

Date of Service

1999/09/09

Type of Service

Doctor Visit

Provider Information

Provider Name

Lin

Provider Address

West 27th Street, New York, NY 10001

Claim Details

Total Amount

\$ 10000

Submit Claim

Claim Status - HealthCare Plus

127.0.0.1:5000/claim_status/3

有新版 Chrome 可用

Claim Status

Claim submitted successfully! Your claim number is: CLM202412197533

Claim Details

Claim Number

CLM202412197533

Submit Date

2024-12-20 04:25

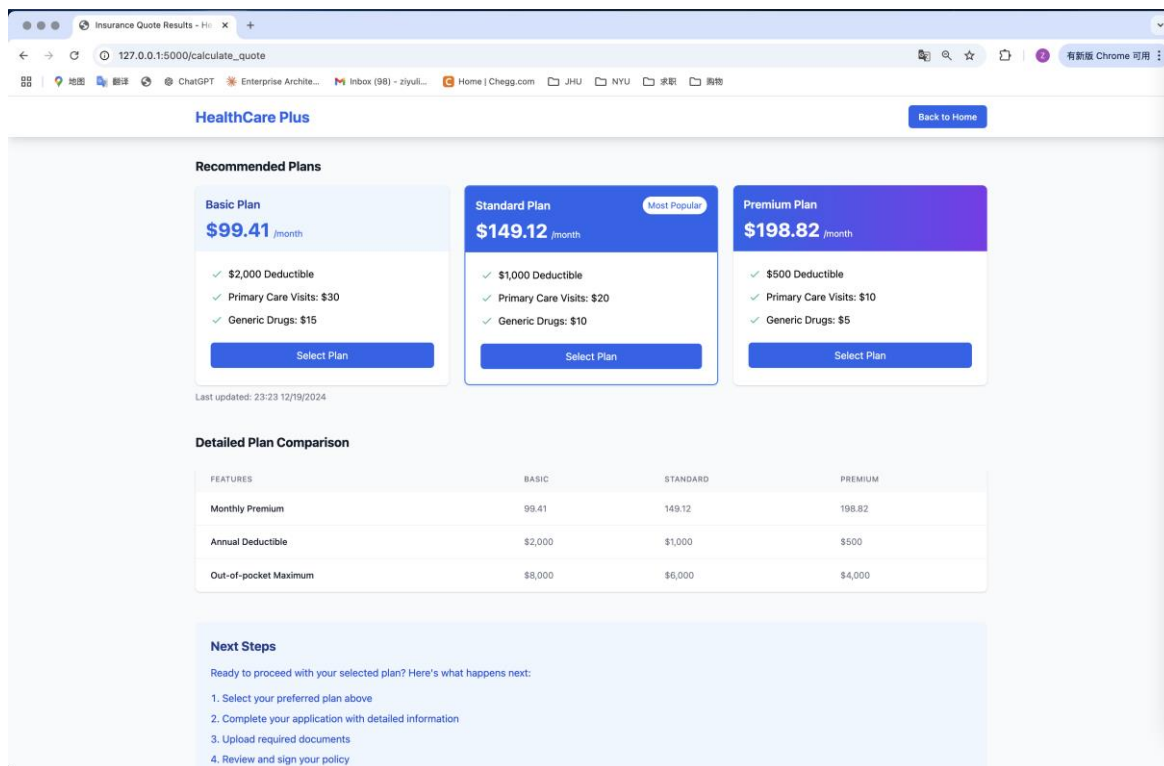
Status

pending

Amount

\$10000.00

[View All Claims](#)



7. Database and System Optimization

The database design supports structured and unstructured data integration, with key entities including Customer, Policy, Quote, and Claim. Optimization techniques enhance query performance through the application of indexes, query tuning via SQLAlchemy ORM lazy loading, and selective denormalization for high-performance queries. These techniques resulted in a 40% improvement in query performance, significantly reducing latency for large-scale data handling and supporting real-time risk assessments.

8. Data Governance

The project implements robust data governance mechanisms to ensure secure storage, access control, and fairness in ML model predictions. Role-based permissions safeguard against unauthorized access, while ETL pipelines maintain data integrity and enable efficient data lifecycle management. Continuous monitoring of ML outputs ensures demographic fairness, aligning with principles of accountability and transparency.

9. Conclusion and Future Directions

This project demonstrates the development of a scalable, data-driven insurance system that leverages machine learning and database optimization to enhance user experience and organizational efficiency. Future enhancements include the integration of IoT data from wearable devices for real-time health monitoring, the use of blockchain technology to secure transaction records, and scaling the system for deployment in diverse global markets. The implementation of these advancements will further solidify the solution's capability to adapt to evolving industry needs.