

HealthAI: Intelligent Healthcare Assistant Using IBM Granite

1.Introduction:

- **Project title :** HealthAI
- **Team leader:** PAVITHRA D
- **Team member :** POOJA V
- **Team member :** SHALINI R

2. Project Overview:

Purpose:

The purpose of HealthAI is to democratize healthcare access by providing instant, AI-driven medical insights to patients anywhere, anytime. With IBM Granite LLM powering its backend, the system is capable of understanding medical prompts and generating contextual responses.

HealthAI focuses on two critical functions as implemented in the current version of the code:

1. **Disease Prediction** – Based on symptoms, it identifies possible conditions and provides general medication suggestions.
2. **Treatment Plan Generator** – Provides personalized treatment advice tailored to the patient's demographics and history.

Goals of the Project:

- Assist patients in understanding their symptoms before consulting a doctor.
- Support healthcare providers with AI-generated treatment ideas.
- Reduce dependency on long hospital visits for basic medical guidance.
- Enhance patient awareness of home remedies and preventive care.

Features:

1. **Conversational Interface** ◦ Function: Natural health Q&A using plain language.

- Example: Patient asks, “What could be the cause of persistent cough?” and receives a possible list of conditions.

2. **Disease Prediction**

- Function: Symptom analysis leading to possible medical conditions. ◦ Example: User enters “fever, cough, fatigue”, and AI suggests possible flu or viral infection with cautionary notes.

3. **Treatment Plan Generator**

- Function: Creates personalized treatment suggestions based on patient details.
- Example: Input: Diabetes, Age 50, Male, History of hypertension. ◦ Output: Lifestyle suggestions, dietary guidelines, and standard treatment options.

3. Architecture

The architecture of HealthAI is structured around simplicity, scalability, and modularity.

Frontend (Gradio):

- Provides an interactive UI with two tabs:
 - *Disease Prediction Tab* – accepts symptom input and generates predictions.
 - *Treatment Plan Tab* – accepts condition, age, gender, and history to create a treatment plan.
- Output displayed in clear, multi-line text boxes.
- Easy integration with additional features (feedback, reports, KPI monitoring).

Backend (Python):

- Written in Python with Hugging Face Transformers and PyTorch.
- Includes modular functions for disease prediction and treatment plan generation.

4. Setup Instructions Prerequisites:

- Python 3.9+
- pip package manager
- Internet access (required to download IBM Granite model)
- GPU (optional, for faster response generation) **Installation**

Process:

1. Install dependencies:

“pip install transformers accelerate gradio torch”

2. Clone repository and open app.py.

3. Run application:

```
python app.py
```

4. Access via Gradio local server:

- Local URL: `http://127.0.0.1:7860`
- Public URL (Colab): Shared link generated by `app.launch(share=True)`

5. Folder Structure:

1. `health_ai.py`

- The main script of the project.
- Contains model initialization, disease prediction function, treatment plan function, and Gradio UI in one file.

2. `requirements.txt`

- Lists dependencies: torch, transformers, accelerate, and gradio.
- Ensures consistent environment setup.

3. `README.md`

- Provides an overview of the project.
- Explains installation steps, how to run the app, and usage guidelines.

6. Running the Application Steps:

1. Start the script in Google Colab or local environment.
2. Load IBM Granite model automatically.

3. Select Disease Prediction Tab → Enter symptoms → Get conditions.
4. Select Treatment Plans Tab → Enter details → Get treatment suggestions.

Example Use Case:

- Input: Symptoms → “headache, nausea, blurred vision”
- Output: Possible conditions include migraine, dehydration, or high blood pressure. Please consult a doctor.

7. API Documentation

- POST /predict-disease – Accepts user-input symptoms in JSON format and returns a list of possible conditions along with recommendations.
- POST /treatment-plan – Accepts patient details such as condition, age, gender, and medical history, and responds with a personalized treatment plan.
- POST /upload-history – Allows uploading past health records to enhance the accuracy of recommendations.

8. Authentication

The current system is open for demo use. For real-world deployment:

- JWT Authentication for user access.
- Role-Based Access Control (RBAC):

- Patient – symptom entry, treatment viewing.
- Doctor – view reports, override AI output.
- Admin – manage application and records.

9. User Interface

- Tabbed Navigation: Disease Prediction & Treatment Plan.
- Textbox Inputs: For symptoms, condition, and medical history.
- Multi-line Outputs: For clear AI responses.
- Accessibility: Designed for ease of use by non-technical patients.

10. Testing Phases of Testing:

1. **Unit Testing:** ○ Checked response formatting, prompt handling.
2. **Manual Testing:** ○ Symptom entry tested with realworld conditions.
3. **API Testing (Future):**
 - Swagger/Postman for endpoint validation.
4. **Edge Case Handling:**
 - Empty input → AI returns “Please enter valid symptoms.”
 - Long inputs → Handled with tokenizer truncation.

11. Output

Coding

3:15



github.com/pavithra-c



25



main - IBM-Project / healthai.py

↑ Top

Code Blame

Raw Download Edit View

```

import gradio as gr
import torch
from transformers import AutoTokenizer, AutoModelForCausalLM

# Load model
model_name = "ibm-granite/granite-3.2-2b-instruct"
tokenizer = AutoTokenizer.from_pretrained(model_name)
model = AutoModelForCausalLM.from_pretrained(
    model_name,
    torch_dtype=torch.float16 if torch.cuda.is_available() else torch.float32,
    device_map="auto" if torch.cuda.is_available() else None
)

if tokenizer.pad_token is None:
    tokenizer.pad_token = tokenizer.eos_token

def generate_response(prompt, max_length=1024):
    inputs = tokenizer(prompt, return_tensors="pt", truncation=True, max_length=512)
    if torch.cuda.is_available():
        inputs = {k: v.to(model.device) for k, v in inputs.items()}

    with torch.no_grad():
        outputs = model.generate(
            **inputs,
            max_length=max_length,
            do_sample=True,
            temperature=0.7,
            pad_token_id=tokenizer.eos_token_id
        )
    response = tokenizer.decode(outputs[0], skip_special_tokens=True)
    return response.strip()

# Disease prediction function
def disease_prediction(symptoms):
    prompt = f"Based on the following symptoms: {symptoms}, suggest possible medical conditions and general treatment suggestion"
    return generate_response(prompt, max_length=1200)

# Treatment plan function
def treatment_plan(condition, age, gender, medical_history):
    prompt = f"Generate personalized treatment suggestions for the following patient information. Include home remedies and general advice"
    return generate_response(prompt, max_length=1200)

# Create Gradio interface
with gr.Blocks() as app:
    gr.Markdown("# 🏥 Medical AI Assistant")
    gr.Markdown("Disclaimer: This is for informational purposes only. Always consult healthcare professionals for medical advice.")

    with gr.Tabs():
        with gr.Tab("Disease Prediction"):
            with gr.Row():
                with gr.Column():
                    symptoms_input = gr.Textbox(
                        label="Enter Symptoms",
                        placeholder="e.g., fever, headache, cough, fatigue...",
                        lines=4
                    )
                    predict_btn = gr.Button("Analyze Symptoms")
                with gr.Column():
                    prediction_output = gr.Textbox(label="Possible Conditions & Recommendations", lines=20)

            predict_btn.click(disease_prediction, inputs=symptoms_input, outputs=prediction_output)

        with gr.Tab("Treatment Plans"):
            with gr.Row():
                with gr.Column():
                    condition_input = gr.Textbox(
                        label="Medical Condition",
                        placeholder="e.g., diabetes, hypertension, migraine...",
                        lines=2
                    )
                    age_input = gr.Number(label="Age", value=30)
                    gender_input = gr.Dropdown(
                        choices=["Male", "Female", "Other"],
                        label="Gender"
                    )
                    history_input = gr.Textbox(
                        label="Medical History",
                        placeholder="Previous conditions, allergies, medications or None",
                        lines=3
                    )
                with gr.Column():
                    plan_btn = gr.Button("Generate Treatment Plan")
                    plan_output = gr.Textbox(label="Personalized Treatment Plan", lines=20)

            plan_btn.click(treatment_plan, inputs=[condition_input, age_input, gender_input, history_input], outputs=plan_output)

app.launch(share=True)

```

3:15



github.com/pavithra-c



25



pavithra-d-05 / IBM-Project


[Code](#)
[Issues](#)
[Pull requests](#)
[Actions](#)
[Projects](#)
[Wiki](#)
[Security](#)
[Insights](#)
[Settings](#)

main • IBM-Project / healthai.py

Go to file



pavithra-d-05 Add files via upload

78a4819 · 1 hour ago

98 lines (82 loc) · 3.98 KB

Code

Blame

Raw



```

1 # -*- coding: utf-8 -*-
2 """healthai.ipynb
3
4 Automatically generated by Colab.
5
6 Original file is located at
7 https://colab.research.google.com/drive/1D5V42F1R1KXdc12N0B83nygU5gJ-TWS
8 """
9
10 !pip install transformers torch gradio -q
11
12 import gradio as gr
13 import torch
14 from transformers import AutoTokenizer, AutoModelForCausalLM
15
16 # Load model
17 model_name = "ibm-granite/granite-3.2-2b-instruct"
18 tokenizer = AutoTokenizer.from_pretrained(model_name)
19 model = AutoModelForCausalLM.from_pretrained(
20     model_name,
21     torch_dtype=torch.float16 if torch.cuda.is_available() else torch.float32,
22     device_map="auto" if torch.cuda.is_available() else None
23 )
24
25 if tokenizer.pad_token is None:
26     tokenizer.pad_token = tokenizer.eos_token
27
28 def generate_response(prompt, max_length=1024):
29     inputs = tokenizer(prompt, return_tensors="pt", truncation=True, max_length=512)
30     if torch.cuda.is_available():
31         inputs = {k: v.to(model.device) for k, v in inputs.items()}
32
33     with torch.no_grad():
34         outputs = model.generate(
35             **inputs,
36             max_length=max_length,
37             do_sample=True,
38             temperature=0.7,
39             pad_token_id=tokenizer.eos_token_id
40         )
41     response = tokenizer.decode(outputs[0], skip_special_tokens=True)
42     return response.strip()
43
44 # Disease prediction function
45 def disease_prediction(symptoms):
46     prompt = f"Based on the following symptoms: {symptoms}, suggest possible medical conditions and general treatment suggestions."
47     return generate_response(prompt, max_length=1200)
48
49 # Treatment plan function
50 def treatment_plan(condition, age, gender, medical_history):
51     prompt = f"Generate personalized treatment suggestions for the following patient information. Include home remedies and over-the-counter medications."
52     return generate_response(prompt, max_length=1200)
53
54 # Create Gradio interface
55 with gr.Blocks() as app:
56     gr.Markdown("🏥 Medical AI Assistant")
57     gr.Markdown("⚠️ Disclaimer: This is for informational purposes only. Always consult healthcare professionals for medical advice.")
58
59     with gr.Tabs():
60         with gr.TabItem("Disease Prediction"):
61             with gr.Row():
62                 with gr.Column():
63                     symptoms_input = gr.Textbox(
64                         label="Enter Symptoms",
65                         placeholder="e.g., fever, headache, cough, fatigue...",
66                         lines=4
67                     )
68                     predict_btn = gr.Button("Analyze Symptoms")
69                 with gr.Column():
70                     prediction_output = gr.Textbox(label="Possible Conditions & Recommendations", lines=20)
71
72             predict_btn.click(disease_prediction, inputs=symptoms_input, outputs=prediction_output)
73
74         with gr.TabItem("Treatment Plans"):
75             with gr.Row():
76                 with gr.Column():
77                     condition_input = gr.Textbox(
78                         label="Medical Condition",
79                         placeholder="e.g., diabetes, hypertension, migraine..."

```


Implementation:

The top screenshot shows the code for generating a treatment plan and the progress of downloading model files. The code includes a function to generate a treatment plan and a button to click it. The progress bar shows the following files and their download status:

File	Progress	Size	Download Speed
tokenizer_config.json	100%	8.88K	0.00-00:00, 18.1MB/s
vocab.json	100%	777K	0.00-00:00, 5.32MB/s
merges.txt	100%	442K	0.00-00:00, 4.89MB/s
tokenizer.json	100%	3.43M	0.00-00:00, 13.9MB/s
added_tokens.json	100%	67.0K	0.00-00:00, 1.64MB/s
special_tokens_map.json	100%	701.7K	0.00-00:00, 13.9MB/s
config.json	100%	786.7K	0.00-00:00, 19.1MB/s
model.safetensors.index.json	29.8K	0.00-00:00, 66.3MB/s	
Fetching 2 files	9%	0.00-00:00, 70%	
model-0002-of-0002.safetensors	100%	67.1M	0.00-00:00, 15.1MB/s
model-0001-of-0002.safetensors	35%	1.76G	0.00-00:00, 23.4MB/s

The bottom screenshot shows the same progress bar further along, followed by a disclaimer and a user interface for the Medical AI Assistant. The progress bar shows the following files and their download status:

File	Progress	Size	Download Speed
tokenizer_config.json	100%	8.88K	0.00-00:00, 18.1MB/s
vocab.json	100%	777K	0.00-00:00, 5.32MB/s
merges.txt	100%	442K	0.00-00:00, 4.89MB/s
tokenizer.json	100%	3.43M	0.00-00:00, 13.9MB/s
added_tokens.json	100%	67.0K	0.00-00:00, 1.64MB/s
special_tokens_map.json	100%	701.7K	0.00-00:00, 13.9MB/s
config.json	100%	786.7K	0.00-00:00, 19.1MB/s
model.safetensors.index.json	29.8K	0.00-00:00, 66.3MB/s	
Fetching 2 files	100%	3.2	0.00-00:00, 124.16MB/s
model-0002-of-0002.safetensors	100%	67.1M	0.00-00:00, 15.1MB/s
model-0001-of-0002.safetensors	100%	5.00G	0.00-00:00, 81.4MB/s
Loading checkpoint shards	100%	2/2	0.00-00:00, 9.31MB/s
generation_config.json	100%	137.1K	0.00-00:00, 15.9MB/s

Colab notebook detected. To show errors in colab notebook, set `debug=True` in `launch()`
* Running on public URL: <https://628322b2f85328578.gradio.live>
This share link expires in 1 week. For free permanent hosting and GPU upgrades, run "gradio deploy" from the terminal in the working directory to deploy to Hugging Face Spaces (<https://huggingface.co/spaces>)

Medical AI Assistant
Disclaimer: This is for informational purposes only. Always consult healthcare professionals for medical advice.

Disease Prediction **Treatment Plans**

Enter Symptoms
e.g., fever, headache, cough, fatigue...

Possible Conditions & Recommendations

Analyze Symptoms

healthai.pytb

https://colab.research.google.com/drive/1D5V42f1R1LXhnci2N0B83nygLSgJ-TMS?authuser=1#scrollTo=c5zR-68KIAe7

healthai.pytb

File Edit View Insert Runtime Tools Help

Q Commands + Code + Text + Run all

generation_config.json: 100% 137/137 [00:00<00:00 15.8kB/s]

Colab notebook detected. To show errors in colab notebook, set debug=True in launch().
* Running on public URL: <https://colab.research.google.com/drive/1D5V42f1R1LXhnci2N0B83nygLSgJ-TMS?authuser=1#scrollTo=c5zR-68KIAe7>

This share link expires in 1 week. For free permanent hosting and GPU upgrades, run "gradle deploy" from the terminal in the working directory to deploy to Hugging Face Spaces (<https://huggingface.co/spaces>)

Medical AI Assistant

Disclaimer: This is for informational purposes only. Always consult healthcare professionals for medical advice.

Disease Prediction

Treatment Plans

Enter Symptoms

e.g., fever, headache, cough, fatigue...

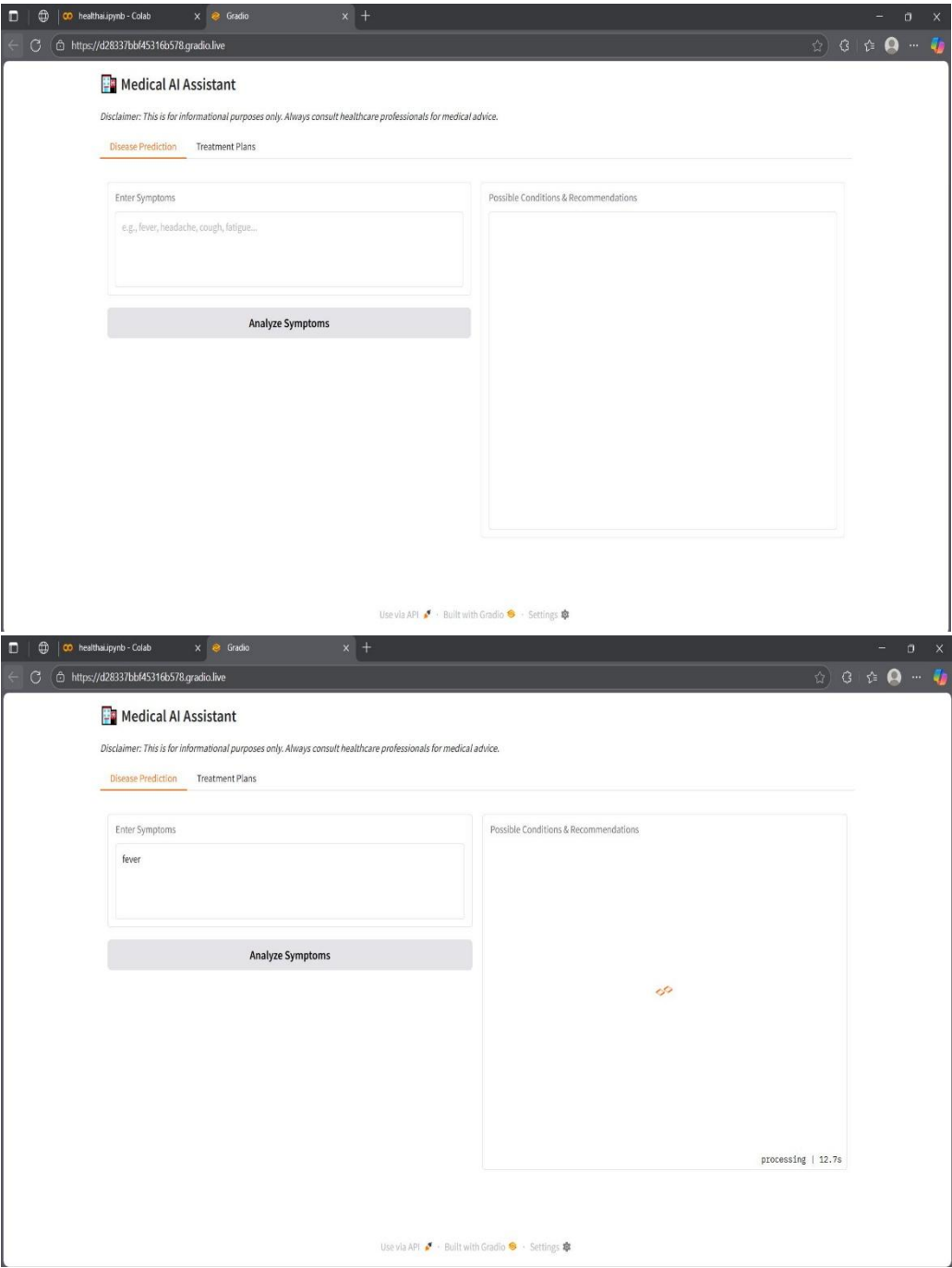
Analyze Symptoms

Possible Conditions & Recommendations

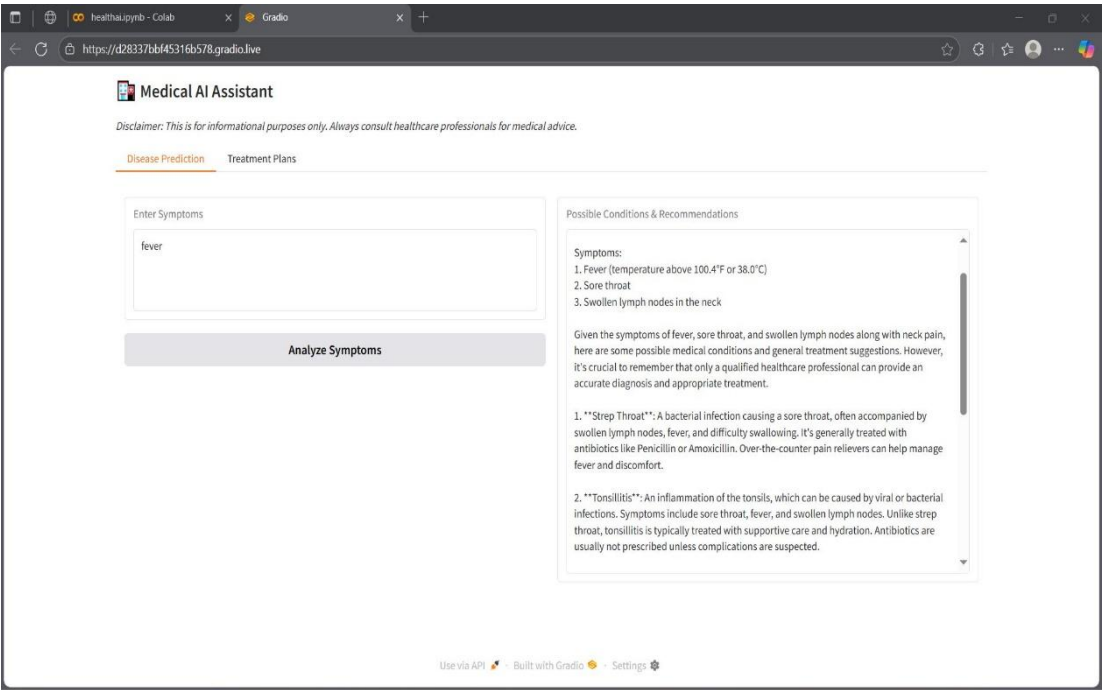
Variables Terminal

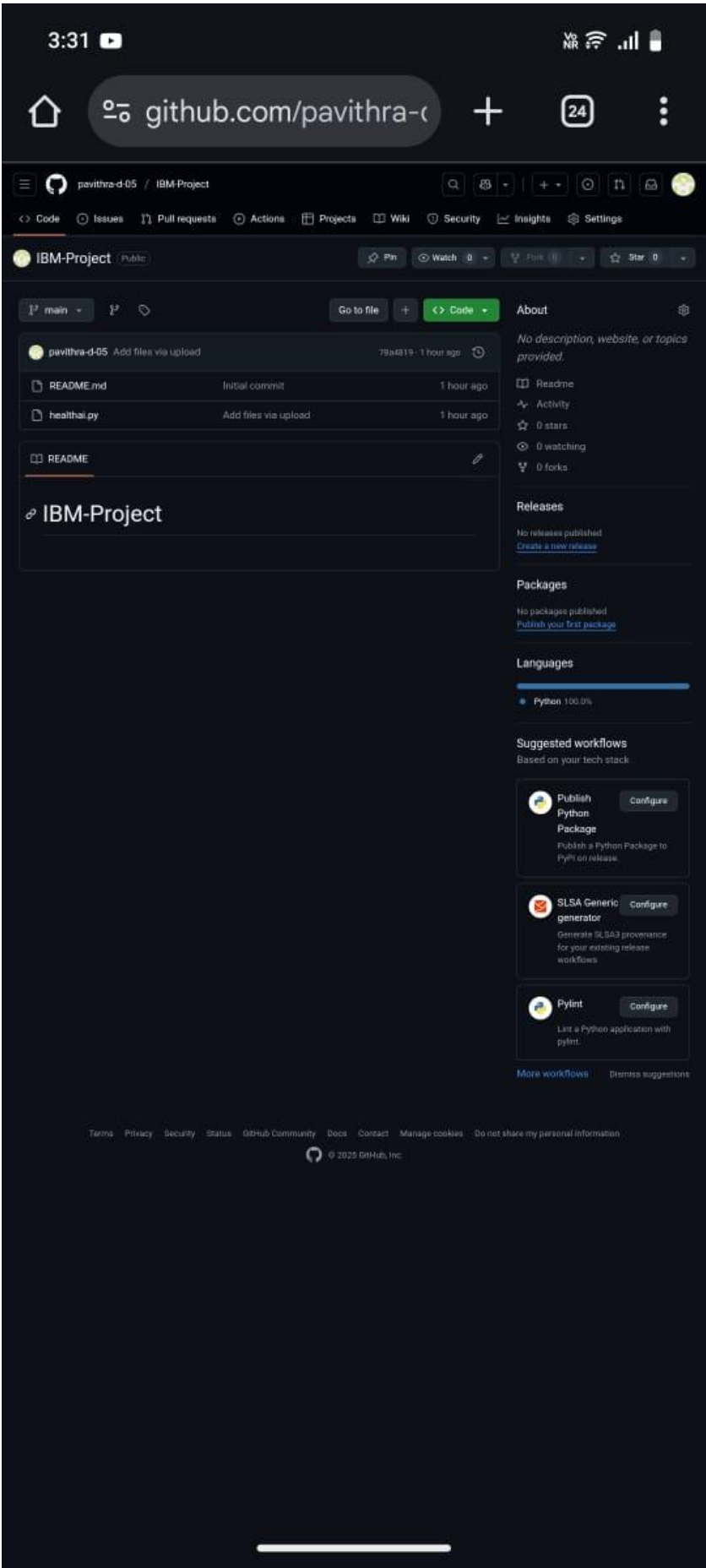
7:45 PM T4 (Python 3)

Output:



final





Treatment plan

12. Known Issues

- Requires stable internet for IBM Granite model.
- Outputs are probabilistic and may vary.
- Does not replace real medical consultation.
- Limited to two major features (prediction & treatment) in current version.

13. Future Enhancements

The current version of HealthAI focuses on symptom-based disease prediction and personalized treatment suggestions. Future improvements based on the existing code framework may include:

- **Medical Report Summarization** – Allow users to upload PDFs or text-based medical reports, which the model can summarize into simple insights.
- **Extended Treatment Plans** – Incorporate additional patient details (diet, lifestyle habits) to generate more comprehensive treatment guidance.
- **Feedback Collection** – Add a feedback tab where patients can share responses, enabling continuous refinement of AI suggestions.

- **Health Data Forecasting** – Introduce forecasting for patient health metrics (e.g., recurring symptoms, predicted recovery timeline).
- **Multilingual Responses** – Enable the model to respond in multiple languages to increase accessibility for nonEnglishspeaking users.