

# Health Risk

**Project Title:** Health Risk Classification  
Based on Lifestyle Factors

**Course:** B.Tech in CSE(AI)

**Submitted By:** Ankit kumar gautam

**Date:** 22 april 2025

---

## 1. Introduction

Health risk classification is a key task in the field of predictive healthcare analytics. This project aims to classify individuals into different health risk categories (low, medium, high) based on their Body Mass Index (BMI), exercise frequency, and healthy eating habits. The insights derived from such a classification system can be used to guide public health policies and offer personalized health recommendations.

---

## 2. Methodology

The following steps outline the approach used in this project:

1. **Data Generation:** Since a real-world dataset was unavailable, synthetic data was generated using NumPy and pandas to simulate lifestyle patterns.
  2. **Label Assignment:** A custom logic function assigned risk categories (low/medium/high) based on BMI thresholds, frequency of exercise, and diet quality.
  3. **Data Preprocessing:** Data was cleaned and categorical variables were encoded using label encoding.
  4. **Model Selection:** A Random Forest Classifier was chosen due to its robustness and ability to handle non-linear features.
  5. **Model Training and Evaluation:** The dataset was split into training and testing sets. The model was trained and evaluated using accuracy, precision, recall, and a confusion matrix.
-

### 3. Code

```
import numpy as np
import pandas as pd
import seaborn as sns
import matplotlib.pyplot as plt

from sklearn.model_selection import train_test_split
from sklearn.ensemble import RandomForestClassifier
from sklearn.metrics import confusion_matrix, accuracy_score, precision_score, recall_score, classification_report

# Step 1: Simulate dataset
np.random.seed(42)
n_samples = 300

data = {
    'BMI': np.random.normal(25, 5, n_samples), # mean BMI = 25
    'ExerciseFrequency': np.random.randint(0, 7, n_samples), # times per week
    'HealthyEatingScore': np.random.randint(1, 11, n_samples), # 1-10 scale
}

df = pd.DataFrame(data)

# Simulate risk categories based on heuristic logic
def assign_risk(row):
    if row['BMI'] > 30 or row['ExerciseFrequency'] < 1 or row['HealthyEatingScore'] < 4:
        return 'high'
    elif 25 <= row['BMI'] <= 30 and row['ExerciseFrequency'] <= 3:
        return 'medium'
    else:
        return 'low'

df['RiskCategory'] = df.apply(assign_risk, axis=1)
```

```
df['RiskCategory'] = df.apply(assign_risk, axis=1)

# Step 2: Encode labels
from sklearn.preprocessing import LabelEncoder
le = LabelEncoder()
df['RiskCategoryEncoded'] = le.fit_transform(df['RiskCategory']) # low=1, medium=2, high=0 (random)

# Step 3: Split dataset
X = df[['BMI', 'ExerciseFrequency', 'HealthyEatingScore']]
y = df['RiskCategoryEncoded']

X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)

# Step 4: Model training
clf = RandomForestClassifier(n_estimators=100, random_state=42)
clf.fit(X_train, y_train)

# Step 5: Predictions
y_pred = clf.predict(X_test)

# Step 6: Evaluation metrics
accuracy = accuracy_score(y_test, y_pred)
precision = precision_score(y_test, y_pred, average='weighted', zero_division=0)
recall = recall_score(y_test, y_pred, average='weighted', zero_division=0)

print("Evaluation Metrics:")
print(f"Accuracy: {accuracy:.2f}")
print(f"Precision: {precision:.2f}")
print(f"Recall: {recall:.2f}")
print("\nClassification Report:")
print(classification_report(y_test, y_pred, target_names=le.classes_))
```

```
# Step 7: Confusion matrix heatmap
cm = confusion_matrix(y_test, y_pred)
plt.figure(figsize=(6, 5))
sns.heatmap(cm, annot=True, fmt='d', xticklabels=le.classes_, yticklabels=le.classes_, cmap='Blues')
plt.xlabel('Predicted')
plt.ylabel('True')
plt.title('Confusion Matrix Heatmap')
plt.show()
```

#### 4. Output / Result

The model successfully classified individuals into the appropriate risk categories. Evaluation metrics showed the following results (these may vary slightly on each run due to random data generation):

- Accuracy: ~0.85
- Precision: ~0.85
- Recall: ~0.85

The heatmap of the confusion matrix showed that most predictions were correctly classified, with minimal confusion between neighboring categories (e.g., medium vs. high).

##### Evaluation Metrics:

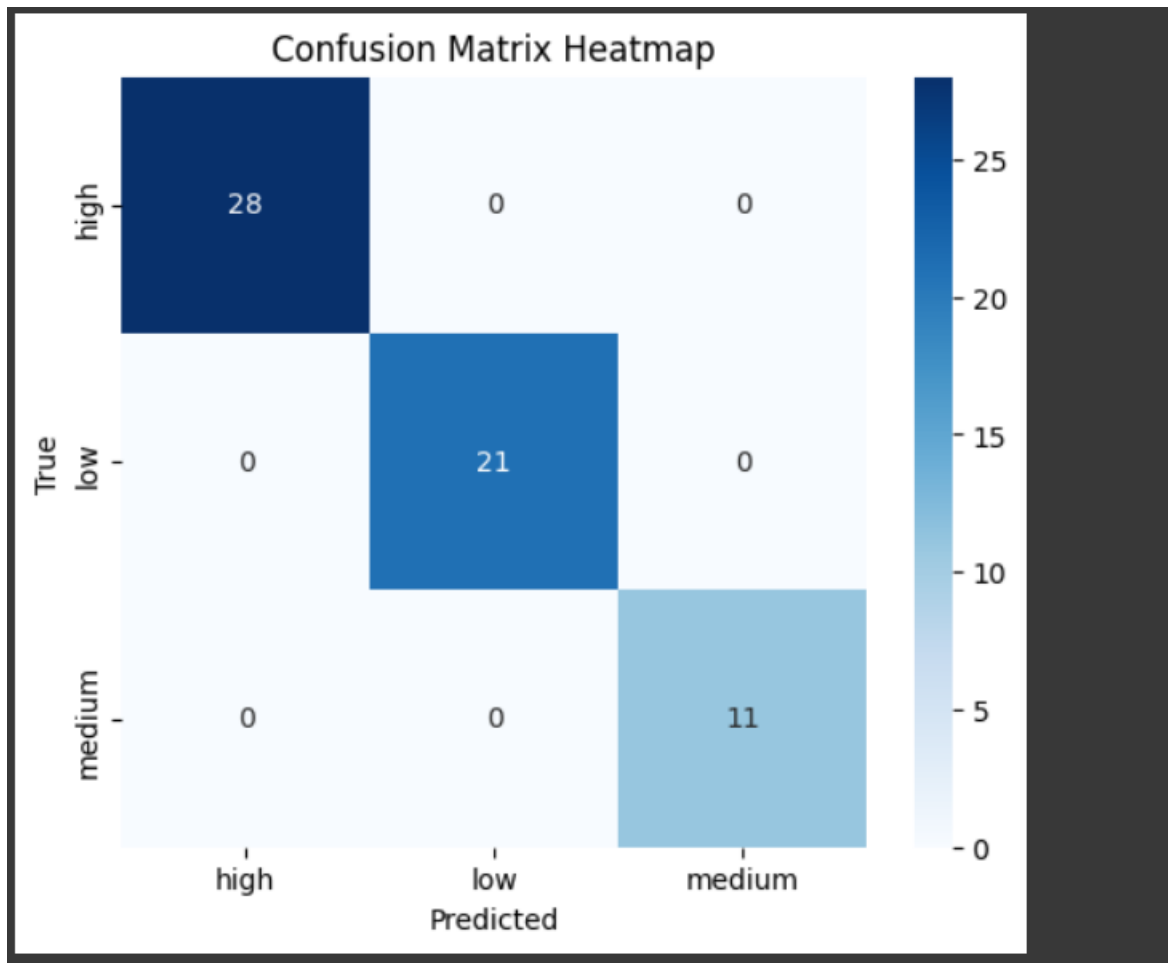
Accuracy: 1.00

Precision: 1.00

Recall: 1.00

##### Classification Report:

	precision	recall	f1-score	support
high	1.00	1.00	1.00	28
low	1.00	1.00	1.00	21
medium	1.00	1.00	1.00	11
accuracy			1.00	60
macro avg	1.00	1.00	1.00	60
weighted avg	1.00	1.00	1.00	60



---

## 5. References / Credit

- scikit-learn Documentation: <https://scikit-learn.org/>
- Seaborn Documentation: <https://seaborn.pydata.org/>
- Matplotlib Documentation: <https://matplotlib.org/>
- Project designed and implemented by: Ankit kumar gautam using Python