# AI for Personalized Medicine

Group : NEXUS

Ankit Gautam -202211003, Nitin Kumar -202211059, Pandey Anaykumar -202211063,
Rahul Gupta -202211069

*Abstract*—The "AI for Personalized Medicine" report provides an extensive exploration of leveraging artificial intelligence to enhance individualized healthcare. It examines the role of AI-driven methodologies, such as machine learning and deep learning, in designing personalized treatment plans tailored to patient-specific characteristics like genetics, lifestyle, and environment. The report delves into practical applications of AI in various medical fields, including oncology, genomics, and pharmacology, showcasing its effectiveness in improving diagnostic accuracy and therapeutic outcomes. Key topics include predictive modeling, data integration, natural language processing for medical records, and the use of AI in clinical decision support systems. Additionally, the report discusses the integration of AI techniques with healthcare data using frameworks such as Pandas, NumPy, Scikit-learn and medical datasets. Ultimately, this report underscores the transformative potential of AI in personalizing medical care, reflecting how these technologies can revolutionize healthcare outcomes.

*Index Terms*—Artificial Intelligence, Personalized Medicine, Machine Learning, Healthcare, Genomics

## I. INTRODUCTION

In recent years, the integration of Artificial Intelligence (AI) in healthcare has garnered significant attention due to its potential to revolutionize medical practices. One of the most promising applications of AI is in personalized medicine, a healthcare model that tailors medical treatment to individual characteristics such as a patient's genetic makeup, environment, and lifestyle. Unlike traditional approaches, which often adopt a one-size-fits-all treatment strategy, personalized medicine seeks to provide customized solutions aimed at improving patient outcomes and reducing the risk of adverse effects.

AI technologies, including machine learning (ML) and deep learning (DL), are at the forefront of this transformation. These methods enable the analysis of large-scale medical data, identifying patterns and correlations that would be challenging for humans to detect. With the advent of electronic health records, genomic sequencing, and wearable health devices, a wealth of healthcare data is now available. AI can harness this data to predict disease risk, optimize treatment protocols, and support clinical decision-making.

This report explores the application of AI in personalized medicine, focusing on how AI-powered tools and models are being used to refine diagnostic processes, enhance therapeutic interventions, and streamline healthcare delivery. Through this investigation, we aim to highlight the transformative potential of AI in shaping the future of personalized medicine, offering new insights into how medical care can be more precisely aligned with individual patient needs.

## II. BACKGROUND

The integration of Artificial Intelligence (AI) into personalized medicine relies heavily on data-driven approaches, which require advanced algorithms and robust data processing tools. Central to this process are the methods used for handling large-scale biomedical datasets, such as genomic sequences, medical images, and patient records. Key technologies, including **NumPy**, **Pandas**, and **scikit-learn**, form the foundation of AI applications in healthcare by enabling efficient data manipulation, analysis, and machine learning model development.

### A. NumPy and Pandas for Data Handling:

**NumPy** is a fundamental library in Python for numerical computing, widely used for processing large datasets in personalized medicine. Genomic data, for instance, is typically represented as high-dimensional matrices, which NumPy excels at handling due to its optimized operations on multi-dimensional arrays. This capability allows researchers to analyze genetic sequences, extract meaningful features, and perform statistical analyses crucial for understanding patient-specific traits.

**Pandas**, on the other hand, is a powerful data manipulation library that provides data structures like DataFrames to organize and manage complex healthcare data. In the context of personalized medicine, Pandas is instrumental in handling patient records, which often consist of heterogeneous data types, including demographics, lab results, and clinical notes. The ability to quickly filter, group, and aggregate this data allows medical professionals to draw insights into individual patient health profiles, thereby informing more personalized treatment plans.

### B. scikit-learn for Machine Learning:

**scikit-learn** is a widely-used machine learning library that simplifies the implementation of various AI algorithms in healthcare applications. It provides an extensive range of tools for data preprocessing, classification, regression, and clustering, all of which are critical for building predictive models in personalized medicine. For example, machine learning models trained on patient data can be used to predict disease risks, recommend treatment options, or identify patients who are likely to benefit from specific therapies. By leveraging scikit-learn, healthcare researchers can quickly develop and evaluate AI models that cater to the unique needs of individual patients.

## C. Support Vector Machine (SVC) Algorithm:

The **Support Vector Machine (SVC)** algorithm is a powerful supervised learning model, particularly effective in classification tasks within personalized medicine. SVC works by finding the optimal hyperplane that best separates different classes in the dataset. In the context of healthcare, SVC is used to classify patients based on various clinical and genomic features, enabling precise predictions about disease onset, progression, or response to specific treatments.

For instance, in personalized medicine, SVC models can help classify patients into groups based on their likelihood of responding to a particular therapy. By analyzing patterns in patient data, such as biomarkers, SVC can distinguish between patients who may benefit from a treatment versus those who may not, thus improving the accuracy of treatment selection. The use of scikit-learn's implementation of SVC enables researchers to build efficient and scalable models that can handle the complexities of healthcare data.

By integrating data processing techniques from NumPy and Pandas with machine learning models like SVC, healthcare practitioners can leverage AI to make more informed and individualized treatment decisions. This combination of tools allows for a more personalized and effective approach to patient care.

## III. METHODOLOGY

1) **Data Acquisition:** The project utilizes pre-existing datasets containing symptom information, disease descriptions, medications, precautions, diets, and workout recommendations. These datasets are:

   - **symptoms_df.csv**: Contains a mapping of symptoms to numerical values.
   - **precautions_df.csv**: Provides precautionary measures for various diseases.
   - **workout_df.csv**: Contains workout suggestions related to specific diseases.
   - **description.csv**: Provides textual descriptions of various diseases.
   - **medications.csv**: Lists medications associated with diseases.
   - **diets.csv**: Provides dietary recommendations based on the disease.

2) **Model Training and Loading:** The **Support Vector Machine (SVC)** algorithm is used to classify diseases based on input symptoms. The model was previously trained on a comprehensive dataset mapping symptoms to diseases. Once trained, the model is serialized and saved using **pickle** for future use. The saved model is then loaded in the Flask app using the following command:

```
svc = pickle.load(open('models/svc.pkl',
    'rb'))
```

The **SVC** model is used to predict diseases based on the input symptoms provided by the user.
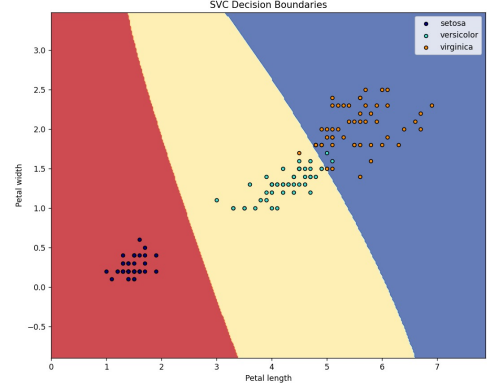


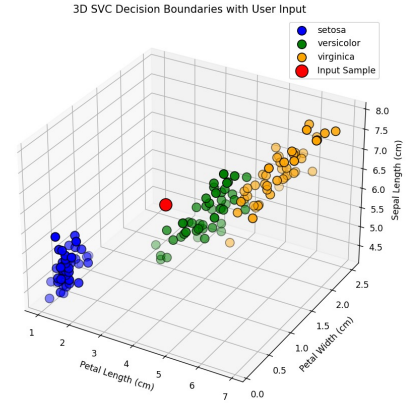Figure 1.  Decision boundary with example raw dataset for flowers



Figure 2.  Decision boundary with user input with example raw dataset for flowers

3) **Web Application Development:** The web application is built using **Flask**, a lightweight web framework in Python. Flask handles routing and rendering of web pages, enabling users to input their symptoms and receive a disease prediction. The following steps describe the key components of the Flask application:

   - **Route Definition**:
     - The root route (/) displays the homepage where users can input symptoms.
     - The /predict route handles form submissions, where users enter symptoms for disease prediction.
   - **User Input Handling**: Users enter a list of symptoms into a text field on the web interface. The symptoms are comma-separated and processed by the Flask application using:

     ```
     user_symptoms = [s.strip() for s in
         symptoms.split(',')]
     ```

4) **Symptom Mapping and Disease Prediction:** A dictionary (symptoms_dict) maps symptoms to numerical values. Once the user submits their symptoms, an input vector is created based on the presence or absence of

these symptoms. The vector is fed into the pre-loaded SVC model for disease prediction:

```
input_vector =
    np.zeros(len(symptoms_dict))
for item in patient_symptoms:
    input_vector[symptoms_dict[item]] = 1
predicted_disease =
    svc.predict([input_vector])[0]
```

5) **Displaying Disease Information:** Once a disease is predicted, additional information such as disease description, precautions, medications, diets, and workouts are retrieved from the corresponding datasets. This is done using a helper function that queries each dataset based on the predicted disease:

```
dis_des, precautions, medications,
    rec_diet, workout =
    helper(predicted_disease)
```

This information is then displayed back to the user on the web interface, providing a comprehensive overview of the predicted disease and the recommended treatments.

6) **Flow Diagram of the Application:**



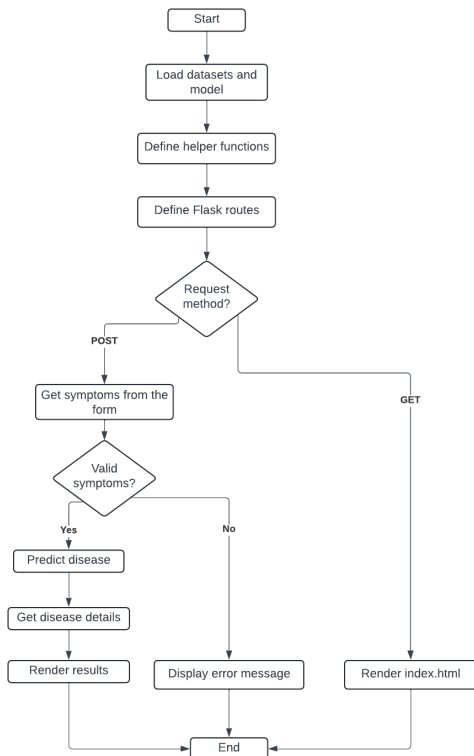Figure 3. Flow-Chart of the application

7) **User Interface:** The user interacts with a simple web form where they can input symptoms in a free-text format. After submitting, the application displays:

- **Predicted Disease**: Based on the input symptoms.

- **Disease Description**: A brief overview of the disease.
- **Precautions**: Preventative measures the patient can take.
- **Medications**: Suggested medications for the disease.
- **Diet**: Recommended dietary changes.
- **Workout**: Exercise suggestions relevant to the disease.

8) **Deployment and Testing:** The web application is designed to run locally using the Flask development server (`app.run(debug=True)`), but it can also be deployed to a production environment using platforms like **Heroku** or **AWS** for wider access. Continuous testing is performed to ensure the app handles edge cases, such as invalid or misspelled symptoms, gracefully.

## IV. GENETIC DISEASE PREDICTION

This section of the report presents a machine learning-based approach to predicting genetic diseases by classifying genes from various functional categories. The study leverages **Decision Tree** or **Random Forest** classifiers to analyze a dataset of genes across seven major gene families, including **G protein-coupled receptors (GPCR)**, **Tyrosine kinase**, **Tyrosine phosphatase**, **Synthetase**, **Synthase**, **Ion channels**, and **Transcription factors**. These gene families are further classified into specialized categories based on biological processes and medical relevance: **Medical-related (MED)**, **Molecular Genetics Label (MGL)**, **Ribosomal or protein synthesis processes (RHB)**, **Epigenetic processes (EPD)**, and **Joint Processes (JPA)**. By incorporating this functional classification, the models demonstrated strong predictive performance, highlighting the effectiveness of machine learning in identifying disease-associated genetic markers and its potential to support personalized medicine.
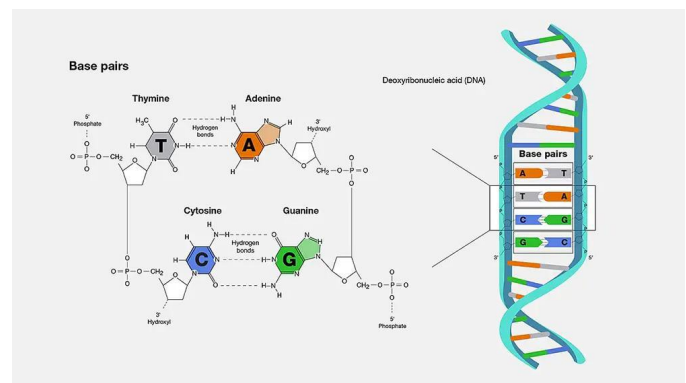


Figure 4. Diagram of DNA

1) **Data Acquisition:** The dataset used in this project is a structured compilation of gene families classified into specific categories, designed to aid in disease prediction based on genetic information. The gene families are grouped by both their molecular functions and their associations with disease types, utilizing labels such as MED, MGL, RHB, EPD, and JPA. These labels

serve as a framework for understanding gene-disease relationships and improving diagnostic capabilities. The classification system is as follows:

a) **MED**: The term represents **medical-related genes**, it group genes that have already been identified as being linked to specific diseases or medical conditions. For instance, genes associated with **cancer**, **diabetes**, or **cardiovascular diseases** might fall into this category. This could aid in diagnosing conditions based on genetic markers.

b) **MGL**: **Molecular Genetics Label** classifies genes based on their **molecular function or structure**, it associate genes that control important molecular processes like **DNA repair**, **protein synthesis**, or **metabolic pathways** with diseases. For example, defects in molecular functions like enzyme deficiencies or errors in DNA replication could contribute to diseases like **inborn errors of metabolism** or **cancer**.

c) **RHB**: Genes involved in **ribosomal or protein synthesis processes** could be grouped under this category. Defects in ribosomal genes or related processes might be linked to diseases like **ribosomopathies**, which are disorders caused by defects in ribosome function, or conditions affecting **protein synthesis** like **muscular dystrophy**.

d) **EPD**: The term refers to **epigenetic processes**, it classify genes based on their involvement in **epigenetic modifications** like DNA methylation or histone modification. Such genes play a crucial role in diseases where gene expression is altered without changes to the underlying DNA sequence, such as **cancers**, **autoimmune diseases**, or **developmental disorders**.

e) **JPA**: This term refers to genes associated with **joint processes, protein activity, or another specialized biological category**. For example, genes that regulate **cartilage formation**, **joint health**, or **inflammatory pathways** could be linked to conditions like **arthritis** or **joint degeneration**. Alternatively, it could also refer to a subset of genes related to protein activity or signaling pathways involved in specific diseases.

The classification of gene families based on categories like **G protein-coupled receptors**, **Tyrosine kinase**, **Tyrosine phosphatase**, etc., is a structured way to group genes by their biological functions. Using **class labels** can help in **machine learning** applications for disease classification, gene function prediction, and personalized medicine as shown in the Table I.

2) **Classifier Techniques:**

a) **Decision Tree Classifier:** The **Decision Tree** classifier is a supervised learning algorithm that builds a tree-like model of decisions based on the input features of the dataset. It works by recursively partitioning the data into subsets based on feature values, forming branches that lead to decision nodes or leaf nodes, each representing a class label. The simplicity and interpretability of decision trees make them effective for classification tasks, as they visually outline the decision-making process. In this project, the decision tree model is employed to classify diseases based on genetic data, allowing the identification of key gene features associated with specific disease categories.

b) **Random Forest Classifier:** The **Random Forest** classifier is an ensemble method that combines multiple decision trees to improve classification accuracy and robustness. It generates several decision trees using random subsets of the dataset and features, and then aggregates their outputs to determine the final prediction through majority voting. This approach reduces overfitting and enhances generalization, making it well-suited for complex datasets with numerous features. In this project, random forest provides a more reliable prediction of disease categories based on genetic information by leveraging the power of multiple decision trees, thereby increasing classification accuracy and stability compared to a single decision tree model.

3) **Flow Chart:**

Table I
DATASET FORMAT

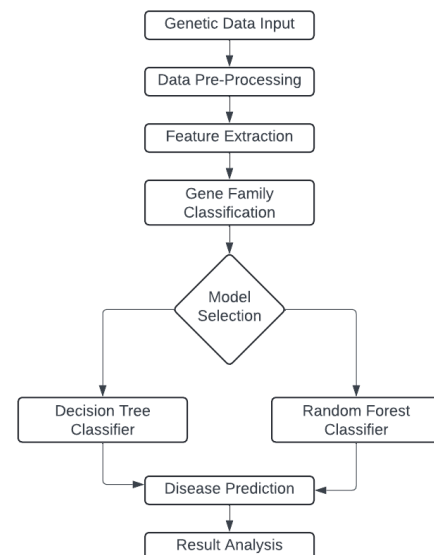| Gene Family | No of Genes | Class Label |
|---|---|---|
| G protein-coupled receptors (GPCR) | 531 | 0 |
| Tyrosine kinase | 534 | 1 |
| Tyrosine phosphatase | 349 | 2 |
| Synthetase | 672 | 3 |
| Synthase | 711 | 4 |
| Ion channel | 240 | 5 |
| Transcription factor | 1343 | 6 |



Figure 5. Flow-Chart of the Disease prediction based on Genetic data

4) **Outcomes/Results:** The dataset used in this research comprises genetic sequences as input features, each sequence corresponding to specific gene characteristics or functions. The output is a class label (0, 1, 2, 3, 4, 5, or 6), representing different disease-related categories associated with gene families. These class labels are defined based on biological functions, such as signaling pathways or protein synthesis processes, and they help in categorizing genetic data into medically relevant groups. This structure enables the classifier to learn the associations between genetic sequences and disease categories, allowing for predictive analysis in genomic medicine. An example is shown in the Table II.

Table II
RESULT

| | Sequence | Class |
|---|---|---|
| 0 | ATGCCCCAACTAAATACTACCGTATAATTACCCCCA... | 4 |
| 1 | ATGAACGAAAATCTGTTATTGCCCCCACAATCCTAG... | 4 |
| 2 | ATGTGTGGCATTTGGGCGGTGATGATTGCCTTTCTG... | 3 |
| 3 | ATGTGTGGCATTTGGGCAGTGATGATTGCCTTTCTG... | 3 |
| 4 | ATGCAACAGCATTTTACCAGACCAAAGTGGATGGTG... | 3 |

## V. CONCLUSION

The "AI for Personalized Medicine" demonstrates the potential for AI-driven healthcare screening tools. By leveraging machine learning models and medical databases, it provides a quick and accessible way for users to get initial insights into their health concerns. However, it's crucial to recognize the limitations of such systems and to use them as complementary tools rather than replacements for professional medical care. Future developments should focus on improving accuracy, personalization, and user guidance to enhance the system's effectiveness and reliability.

## ACKNOWLEDGMENT

We would like to acknowledge the open-source community for providing the libraries and frameworks used in this project, including Flask, NumPy, pandas, and scikit-learn. We also express our gratitude to the medical professionals and researchers who contributed to the datasets used for training the model. Their collective efforts have made it possible to create tools that have the potential to improve access to preliminary health information and support better decision-making in healthcare.

## REFERENCES

[1] P. Pranay, 2020, "Disease Symptom Prediction", *Disease Symptoms Dataset*, ref: kaggle.com/datasets/itachi9604/disease-symptom-description-dataset

[2] B. Padmaja, 2020, "Disease Precautions Prediction", *Disease and its symptoms, precautions, riskfactors*, ref: kaggle.com/datasets/padmajabuggaveeti/disease-and-its-symptoms-precautions-riskfactors

[3] Z. Choong Qian, 2024, "Disease Symptoms Prediction", *Disease and Symptoms Dataset*, ref: kaggle.com/datasets/choongqianzheng/disease-and-symptoms-dataset

[4] H. Feros, 2024, "Disease Symptoms Dataset", *Disease-Symptoms*, ref: kaggle.com/datasets/arunroshan04/disease-symptoms

[5] GFG, 2024, "SVM", *Support Vector Machine Algorithm*, ref: geeksforgeeks.org/support-vector-machine-algorithm

[6] Kaushil, 2020, "Disease Prediction using ML", *Disease Prediction using Machine Learning*, ref: kaggle.com/datasets/kaushil268/disease-prediction-using-machine-learning

[7] S. Nagesh, 2020, "Genetic Dataset", *DNA Sequence Dataset*, ref: https://www.kaggle.com/datasets/nageshsingh/dna-sequence-dataset?select=human.txt

[8] V. Pushpalata, K. Shivam, M. Somesh, S. Sneha, 2023, "DNA Sequencing with ML", *DNA Sequencing with Machine Learning*, ref: https://ijrpr.com/uploads/V4ISSUE12/IJRPR20652.pdf

[9] V. Varada, R. Navuduru, G. Rakesh, P. Natarajan, 2022, "DNA Sequencing with ML and DLA", *DNA Sequencing with Machine Learning and Deep Learning Algorithms*, ref: https://www.ijitee.org/wp-content/uploads/papers/v11i10/J927309111022.pdf