

# GA for Feature Selection in Training Cancer Prediction ML Models

A presentation by  
Soumaya Adabala  
Enrollment Id 24EE62R03  
M.Tech Control Systems Engineering  
IIT Kharagpur

September 15, 2024



# Contents

- 1 Introduction
- 2 Flowchart
- 3 Fitness Function
- 4 Software Demonstration
- 5 Results
- 6 Conclusions
- 7 References

# Introduction

# The Problem

## **“Curse of dimensionality”**

Unnecessary data leads to excessive complications.

- Genetic testing is done to analyze DNA to find mutations linked to cancer risk.
- Bioinformatics data includes diverse features like gene expressions, and protein sequences.
- Extensive testing may lead to higher medical expenses.
- Complex data can create uncertainty about cancer prognosis.

# The problem

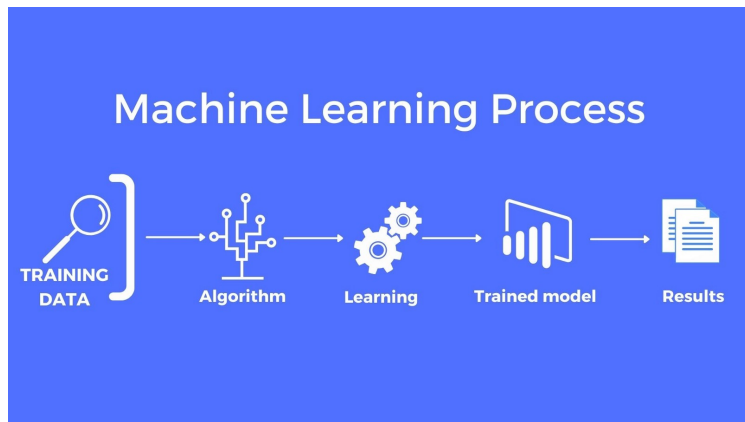


Figure: ML Process

# The Solution

- **Feature Selection:** Selecting the right subset of data.

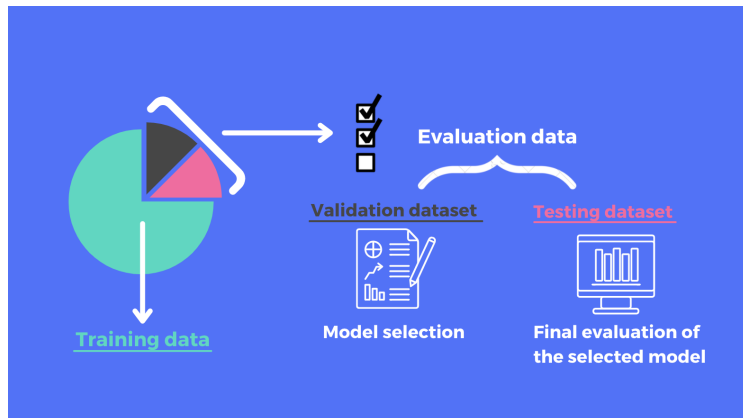


Figure: Data Validation

# The Solution

- Genetic Algorithms (GA) are used in feature selection to optimize the selection of relevant features by evolving subsets that maximize model performance.
- GAs help in reducing dimensionality by identifying the most significant features, thereby improving computational efficiency and model accuracy.
- **Why GA?**
  - flexible
  - adaptable
  - robust
  - efficiently handling complex and nonlinear feature interactions.

# Flowchart



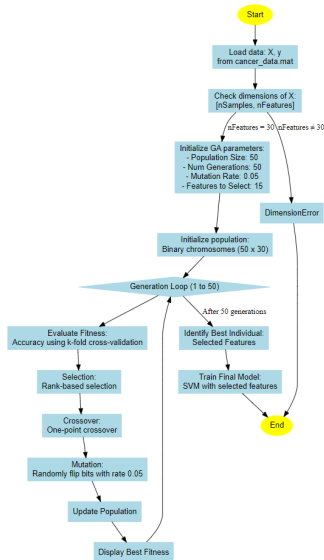


Figure: Flowchart for feature selection using GA

# Fitness Function

# Fitness Function

- **Objective:**
  - Evaluate the performance of feature subsets based on classification accuracy.
- **Inputs:**
  - **Feature Matrix  $\mathbf{X}$ :**  $n \times m$  (samples x features)
  - **Target Labels  $\mathbf{y}$ :**  $n \times 1$
  - **Selected Features  $\mathbf{s}$ :**  $m \times 1$  binary vector (1 if selected, 0 if not)

- **Feature Subset Extraction:**

$$\mathbf{X_s} = \mathbf{X}(:, \mathbf{s})$$

- **Model Training & Evaluation:**

- Use  $k$ -fold cross-validation to evaluate model performance.
- Train a classifier (SVM Support Vector Machine) and compute accuracy for each fold.

- **Accuracy Calculation:**

$$\text{Accuracy}^i = \frac{\sum_{j=1}^{|\mathbf{X}_{\text{test}}^i|} (\text{predicted}_j == \text{actual}_j)}{|\mathbf{X}_{\text{test}}^i|}$$

- **Fitness Score:**

$$\text{Fitness}(\mathbf{s}) = \frac{1}{k} \sum_{i=1}^k \text{Accuracy}^i$$

# Software Demonstration

# Results

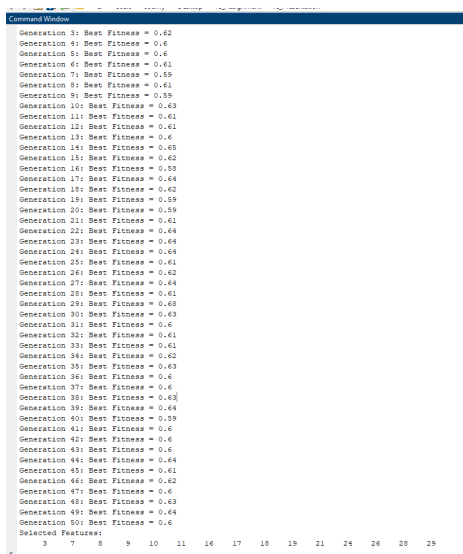


Figure: Results

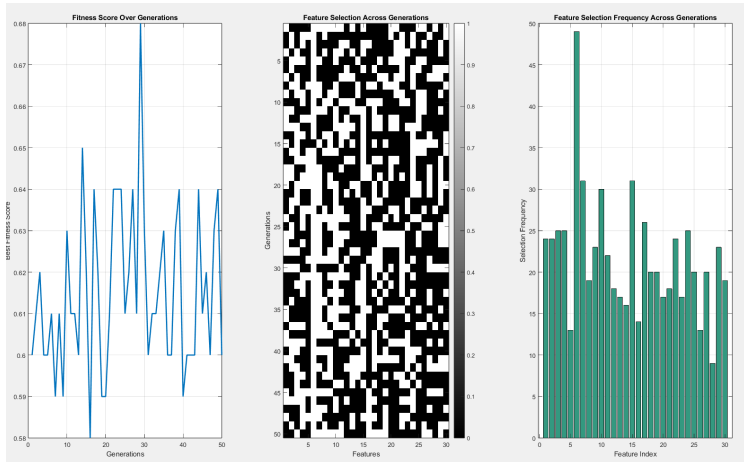


Figure: Results



# Conclusions

# Conclusions

- Genetic Algorithms optimize feature selection by evolving subsets to enhance model performance, reduce dimensionality, and improve both computational efficiency and accuracy.
- Applications include.
  - Healthcare Diagnosis
  - Customer Segmentation
  - Predictive Maintenance

# References

# References

- Liu, B. G., Xu, L. J., Wang, Y. H., & Tang, J. H. (2012). *Genetic Algorithm-Based Feature Selection for Classification: A Comparative Study*. Presented at the IEEE International Conference on Systems, Man, and Cybernetics (SMC). Available at IEEE Xplore.
- García, J. S., Gómez, J. A., & Carrillo, A. J. L. (2013). *Feature Selection Using Genetic Algorithms: A Review*. IEEE Transactions on Evolutionary Computation. Available at IEEE Xplore.
- Pal, S. K., Sinha, S. K., & Chaudhuri, B. B. (2011). *Feature Selection for Classification Using Genetic Algorithms*. Pattern Recognition. Available at ScienceDirect.

Thank You