

## DATA ANALYTICS PROJECT PRESENTATION FEATURE EXTRACTION USING GENETIC ALGORITHM

Shreyas Battula PES1UG20CS405 Kiran H Kademani PES1UG20CS654 Vinay Kumar S PES1UG20CS701



#### PROJECT PRESENTATION

- Feature selection has become a fundamental step of data processing for training a machine learning model due to the increasing amount of high dimensional data.
- There are a huge number of state-of-the-art algorithms that aim to optimize feature selection including genetic algorithms. While not the best performing technique, genetic algorithms provide an exciting solution based on evolution and are widely used in fields such as robotics, marketing or medicine.
- It is worth mentioning that the selection of individuals for reproduction is stochastic, so the best individuals will not always be selected, helping the algorithm not to fall into local minimums.

#### Implementation Genetic Algorithm Steps



Modules Used:

1.Initialise Random Population

2. Fitness Evaluation

3. Selection of Best Parents

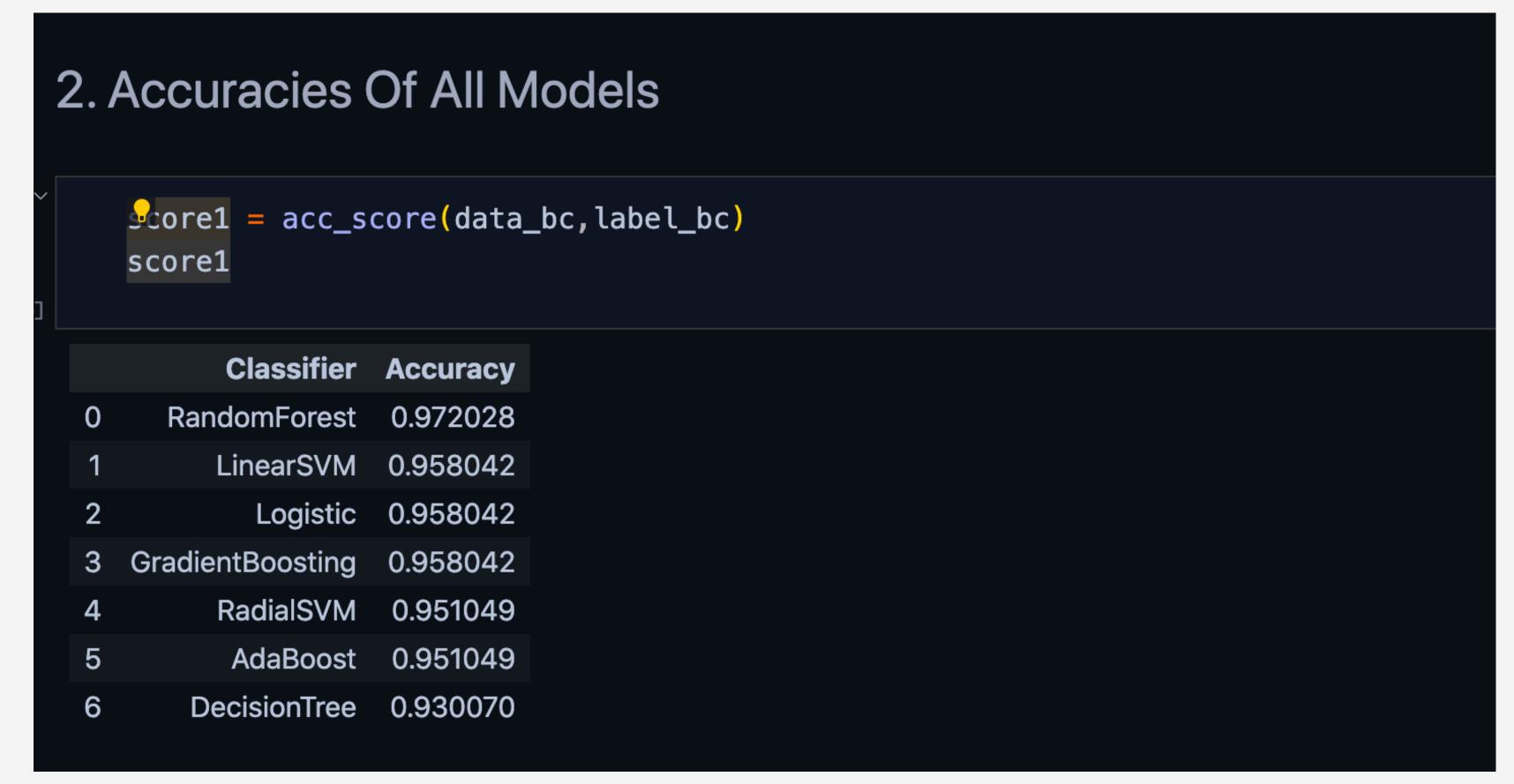
4.CrossOver

5. Mutation

6.Generation

#### OUTPUTS

#### Performance of Models when All Features are Used



### Performance of Random Forest Classifier When GA is used for Features Selection. -An improvement of 2%

#### List of Selected Features For Every Generation

The features selected for generation 1 are

['radius\_mean', 'texture\_mean', 'smoothness\_mean', 'compactness\_mean', 'symmetry\_mean', 'radius\_se', 'perimeter\_se', 'area\_se', 'smoothness\_se', 'concavity\_se', 'concave points\_se', 'symmetry\_se', 'radius\_worst', 'area\_worst', 'smoothness\_worst', 'concavity\_worst', 'symmetry\_worst']

The features selected for generation 2 are

['radius\_mean', 'texture\_mean', 'perimeter\_mean', 'smoothness\_mean', 'compactness\_mean', 'fractal\_dimension\_mean', 'radius\_se', 'texture\_se', 'perimeter\_se', 'area\_se', 'smoothness\_se', 'compactness\_se', 'concavity\_se', 'concave points\_se', 'symmetry\_se', 'fractal\_dimension\_se', 'radius\_worst', 'area\_worst', 'smoothness\_worst']

The features selected for generation 3 are

['radius\_mean', 'texture\_mean', 'smoothness\_mean', 'compactness\_mean', 'fractal\_dimension\_mean', 'radius\_se', 'perimeter\_se', 'area\_se', 'smoothness\_se', 'compactness\_se', 'concavity\_se', 'concave points\_se', 'symmetry\_se', 'fractal\_dimension\_se', 'radius\_worst', 'smoothness\_worst', 'concave points\_worst']

The features selected for generation 4 are

['concavity\_mean', 'texture\_se', 'perimeter\_se', 'area\_se', 'smoothness\_se', 'compactness\_se', 'concave points\_se', 'perimeter\_worst', 'area\_worst', 'smoothness\_worst', 'compactness\_worst', 'concavity\_worst', 'fractal\_dimension\_worst']

The features selected for generation 5 are

['perimeter\_mean', 'concavity\_mean', 'texture\_se', 'perimeter\_se', 'area\_se', 'compactness\_se', 'concavity\_se', 'concave points\_se', 'smoothness\_worst', 'concavity\_worst', 'fractal\_dimension\_worst']

#### Conclusion



Feature extraction is an important data preprocessing step that allows us to improve the accuracy of our results.

Genetic Algorithms are stochastic thus not always falling into the local minimum making it one of the best algorithms to optimise feature selection

#### Future Work

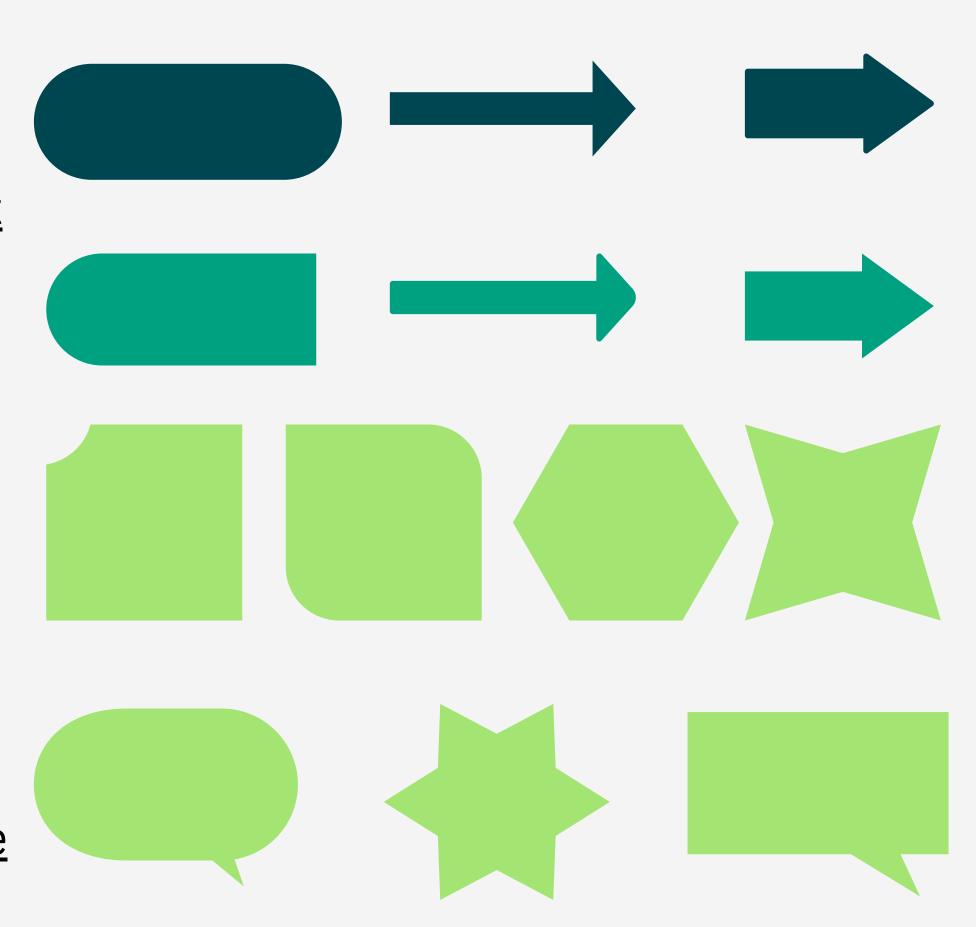
- Compare Model Metrics Against various
   Selection Techniques
- Try to reduce computational time and achieve the required accuracy with minimal number of generations.
- Try to avoid the risk of overfitting of feature selection when observations is insufficient.
- Try A Hybrid GA for larger datasets

#### References

https://towardsdatascience.com/feat ure-selection-using-geneticalgorithms-d3f5fc7bbef1

https://medium.com/analyticsvidhya/feature-selection-usinggenetic-algorithm-20078be41d16

https://www.neuraldesigner.com/blo g/genetic\_algorithms\_for\_feature\_sele ction



# Thank you.