

A stylized scientific illustration on the left side of the slide. It features a large orange circle containing an open book with green and white pages. A magnifying glass with a green handle and frame is positioned over the book. Surrounding the book are various laboratory items: a round-bottom flask with green liquid, a graduated cylinder with green liquid, a thermometer with an orange bulb, and a molecular structure with orange and green spheres connected by lines. The background is a light green color with white and orange abstract shapes.

Classification of Mice on Protein Expression Level

-Team E

Team Leads

Vaibhavi

Lead

Nandita

Co-Lead

Akshar

Co-Lead

Extended Members



Aachal Gupta

Aarjav Jain

Pabitra Goswami

Saili Ashok Dhuri

Pratham Jindal

Ashmit Zanzote

Shubham Arvind Rangari

Abeed Mohammed

Ketan Kumar Sandre

Juweriya Shayhmeeh

Kunal Khumbhar

Sudarshan Sopane

Table of contents



01

Introduction

02

Data set Explanation

03

Steps Explanation

04

Analysis

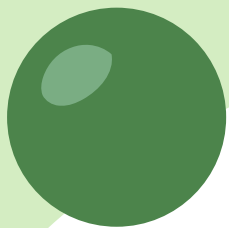
05

Challenges

06

Conclusion





Introduction

Our project, "Classifying Mice Based on Protein Expression Levels," explores the relationship between protein expression, genetics, behavior, and treatments in mice. Using machine learning, we developed a model to classify mice into eight groups based on 77 proteins.

This research provides insights into how Down syndrome affects learning and memory and evaluates the potential benefits of the drug memantine for trisomic mice.



Steps Involved

A

Data
Preprocessing

D

Model Training

B

Exploratory Data
Analysis

E

Model Evaluation

C

Feature Selection

F

Interpretation &
Analysis

Data Preprocessing

An infographic with a light green background. At the top center is the title 'Data Preprocessing' in a bold, dark grey font. Below the title are three white circles, each containing a preprocessing step in bold black text: 'Handling Missing Values' on the left, 'Data Normalisation' in the center, and 'Encoding Categorical Variables' on the right. An open book with yellow pages and green covers is positioned behind the central circle. Various decorative elements are scattered around: a thermometer with orange liquid in the top left, a green sphere in the top right, a green sphere in the bottom right, a green cloud in the top left, a green cloud in the bottom right, a large orange sphere in the bottom left, and a white flask with green liquid in the bottom right. Thin white lines connect some of these elements.

**Handling
Missing
Values**


**Data
Normalisation**

**Encoding
Categorical
Variables**

EDA




Data Distribution




Distribution of protein expression levels across different classes.

Class Balance




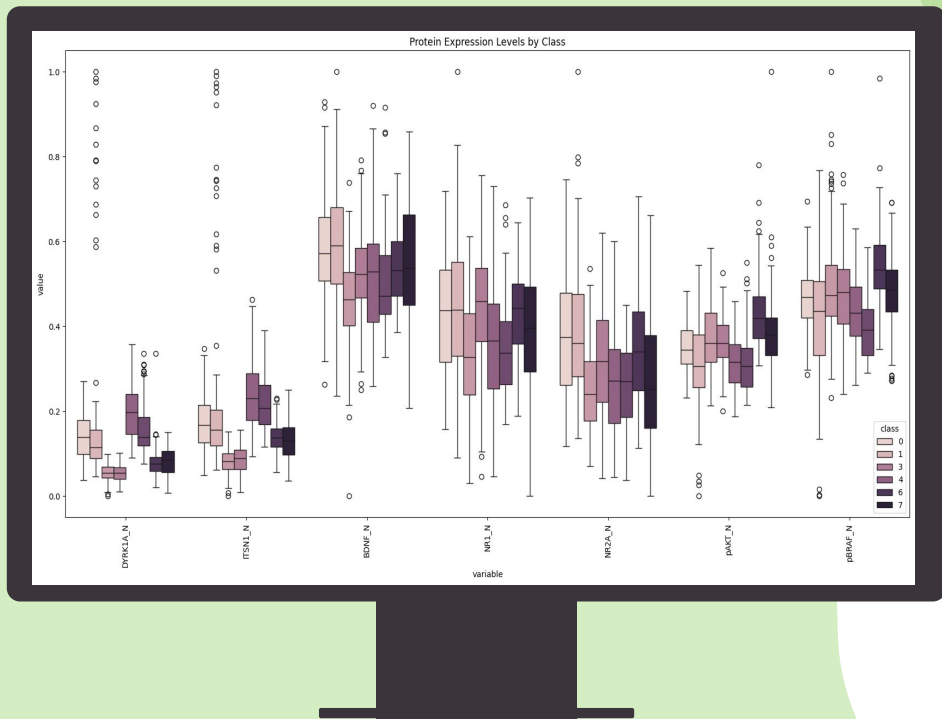
Examination of the distribution of samples across different classes to check for class imbalance.

Missing Values



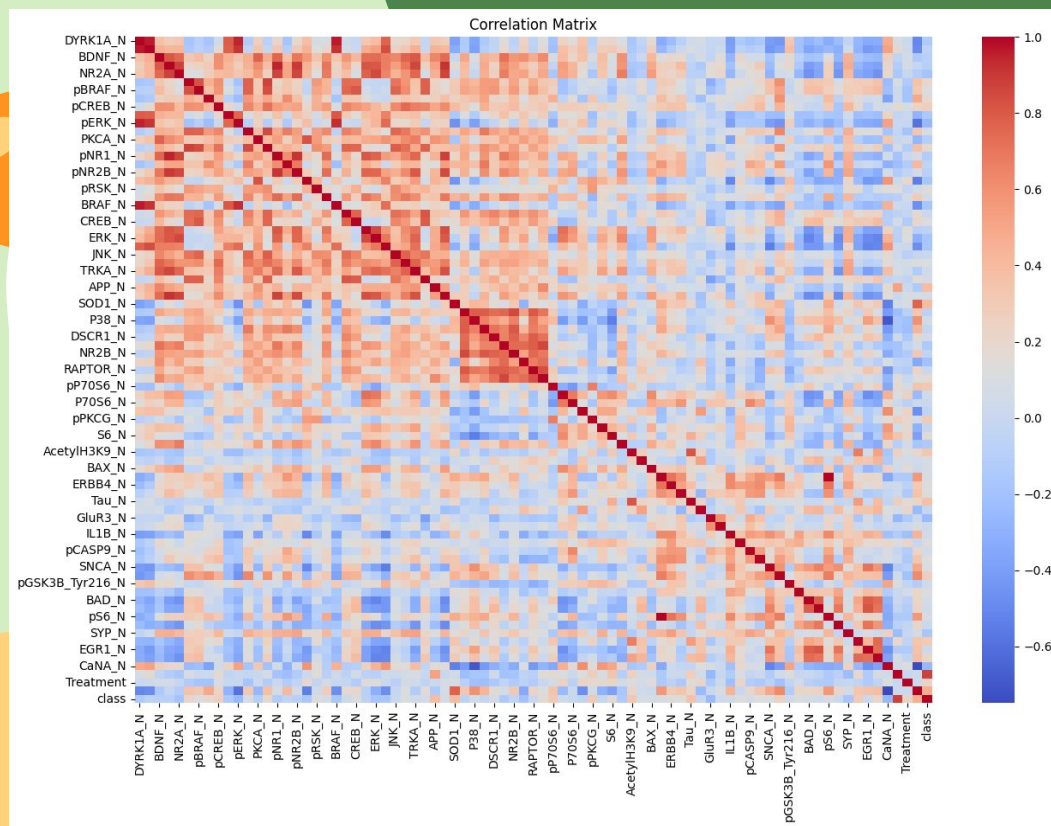
Identification and proportion of missing values in the dataset.



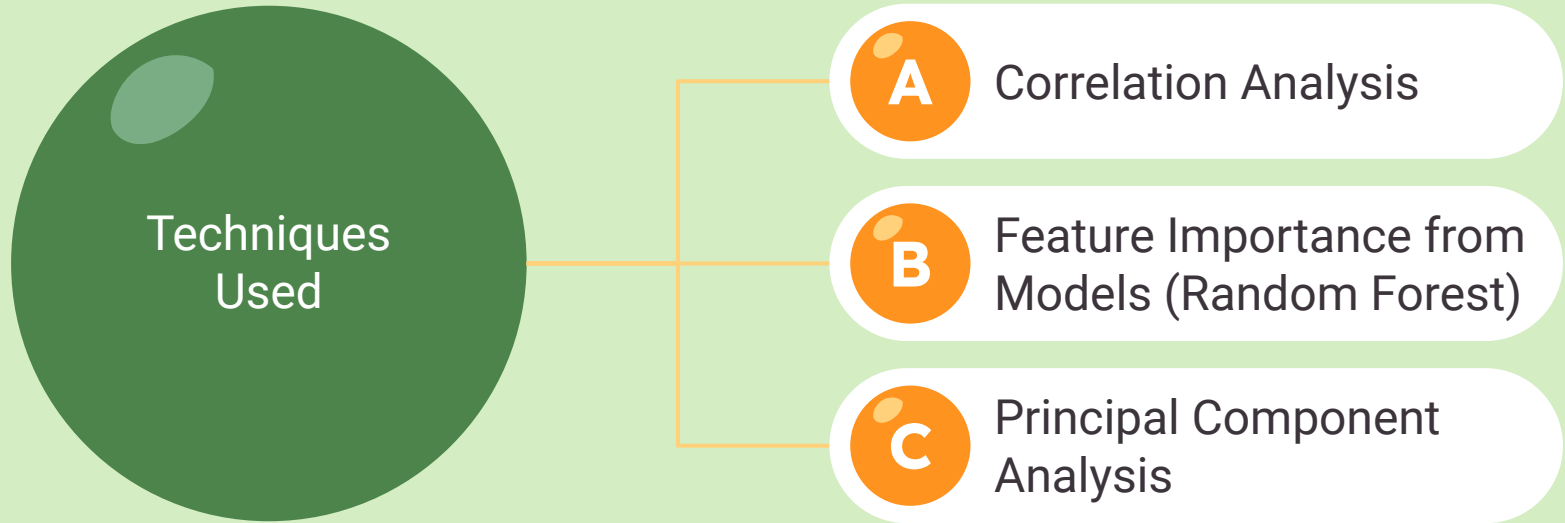


Box Plots

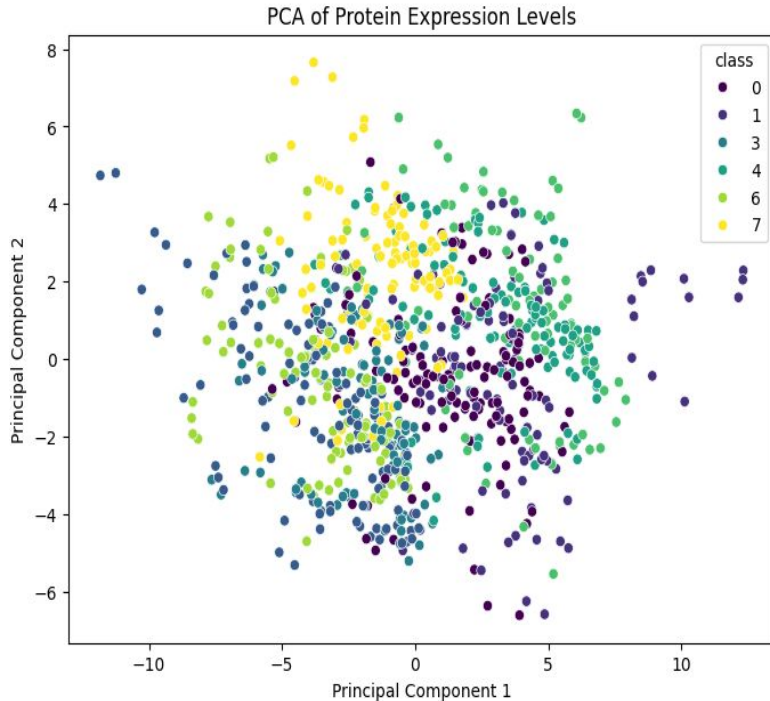
Created to visualise the distribution and variance of protein expression levels across different classes.



Feature Selection

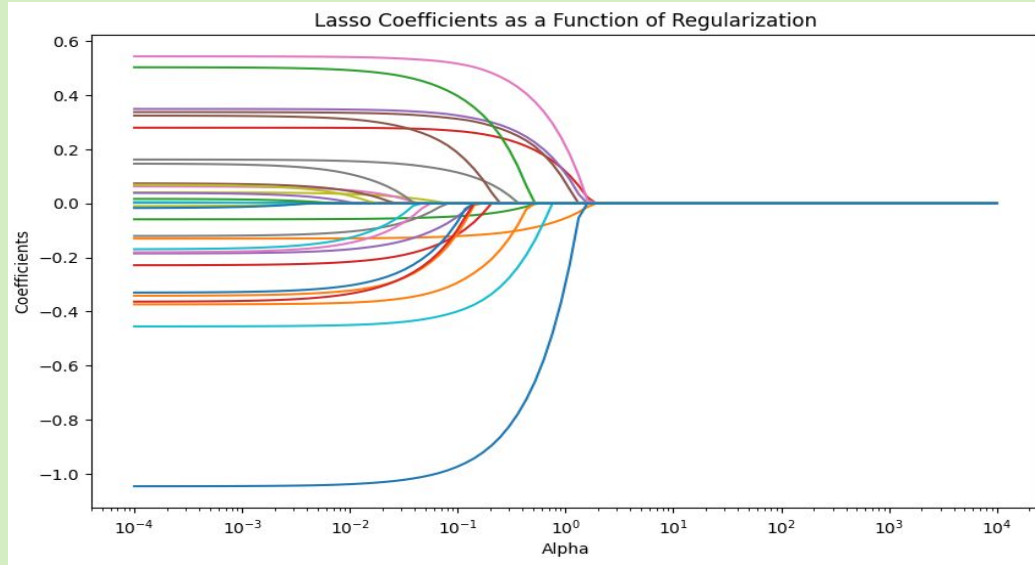


Feature Selection



PCA reduces the complexity of the dataset while retaining most of the variance. Tight, distinct clusters indicate good separability, while overlapping clusters suggest that the classes are not well separated by the chosen features.

Feature Selection



● Use of Lasso

The Lasso coefficients plot shows how different feature coefficients shrink to zero as the regularisation parameter (α) increases, many coefficients shrink to zero, indicating that Lasso regression effectively reduces the number of features by retaining only the most significant ones.

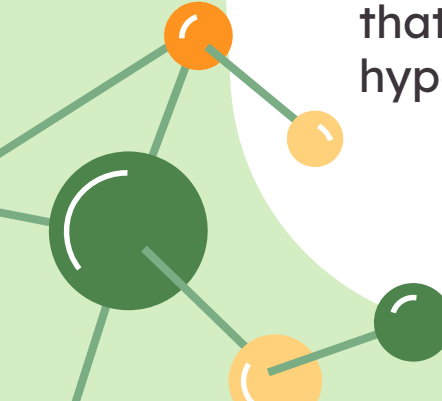
A subset of key proteins was identified as highly discriminative between control and trisomic samples.



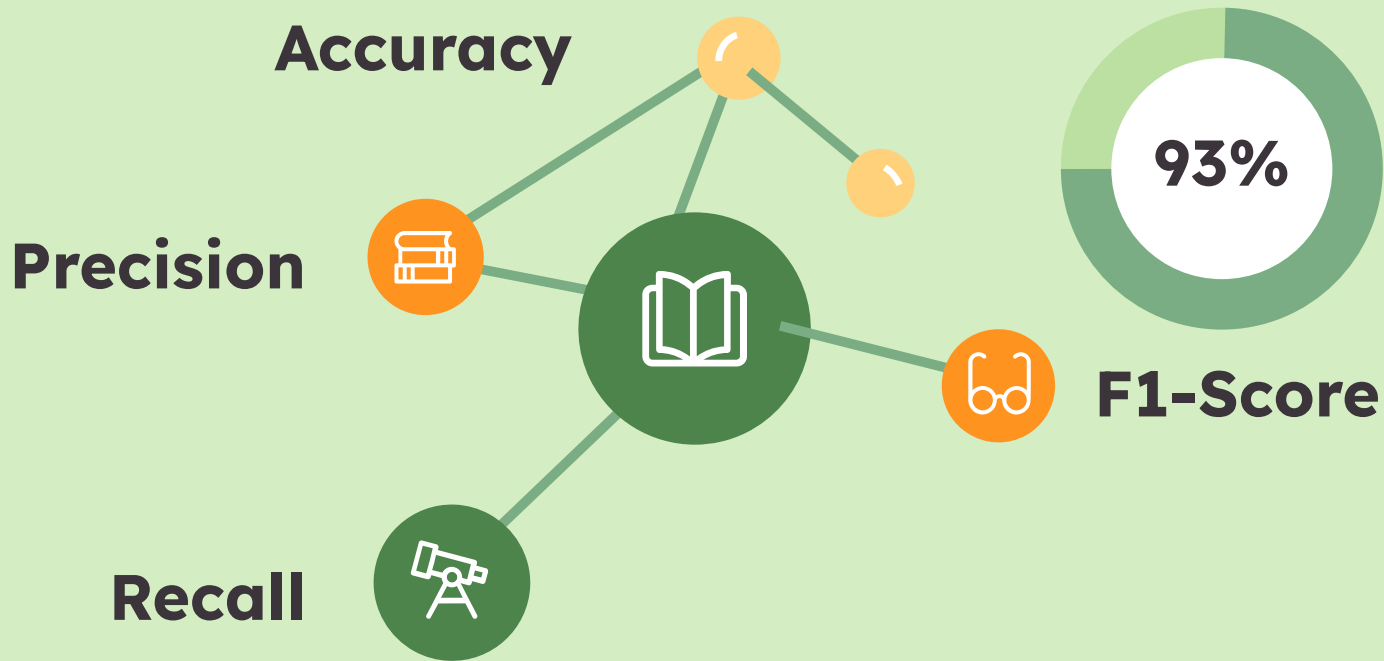
What we used??

Random Forest: An ensemble method that constructs multiple decision trees and outputs the mode of the classes for classification.

HyperParameter Tuning: An exhaustive search method that tests all possible combinations of predefined hyperparameters to find the best-performing model.



Model Evaluation



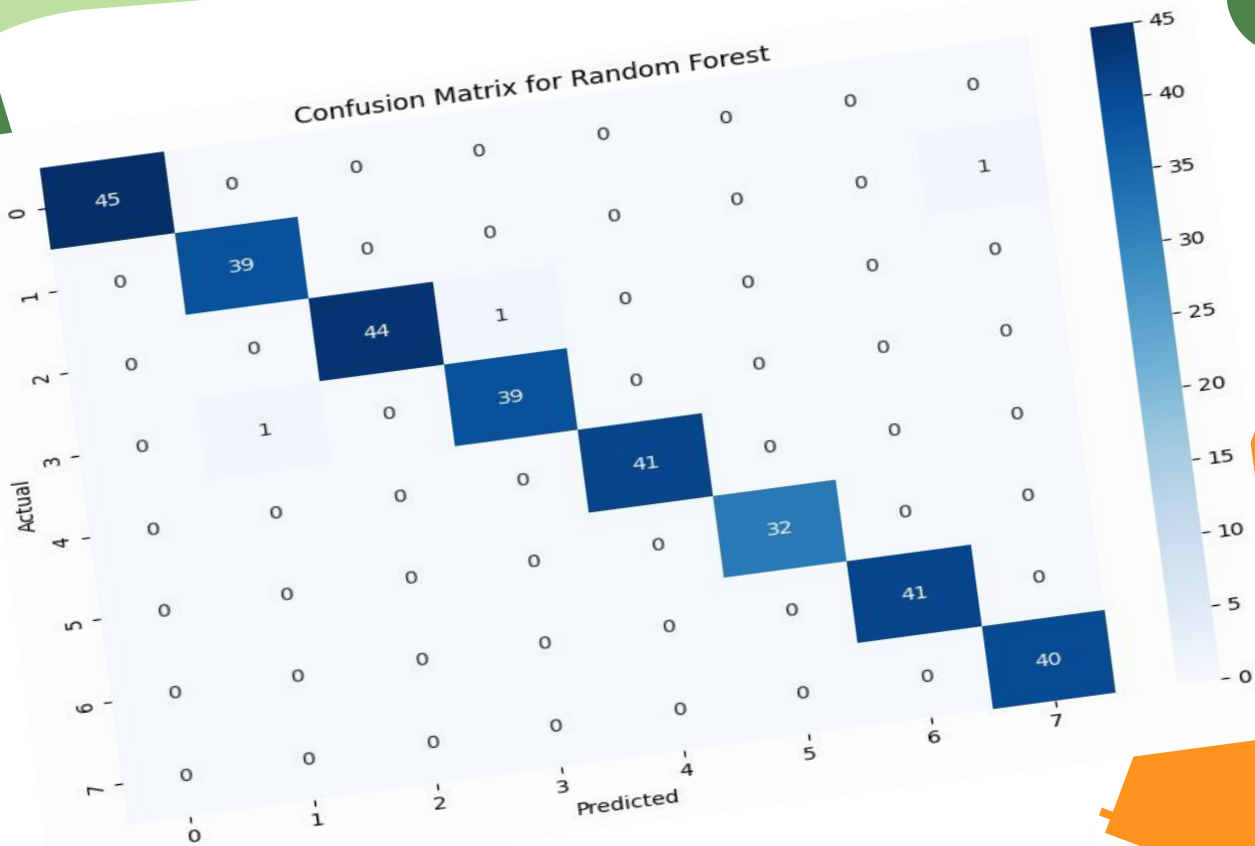
Here's a table

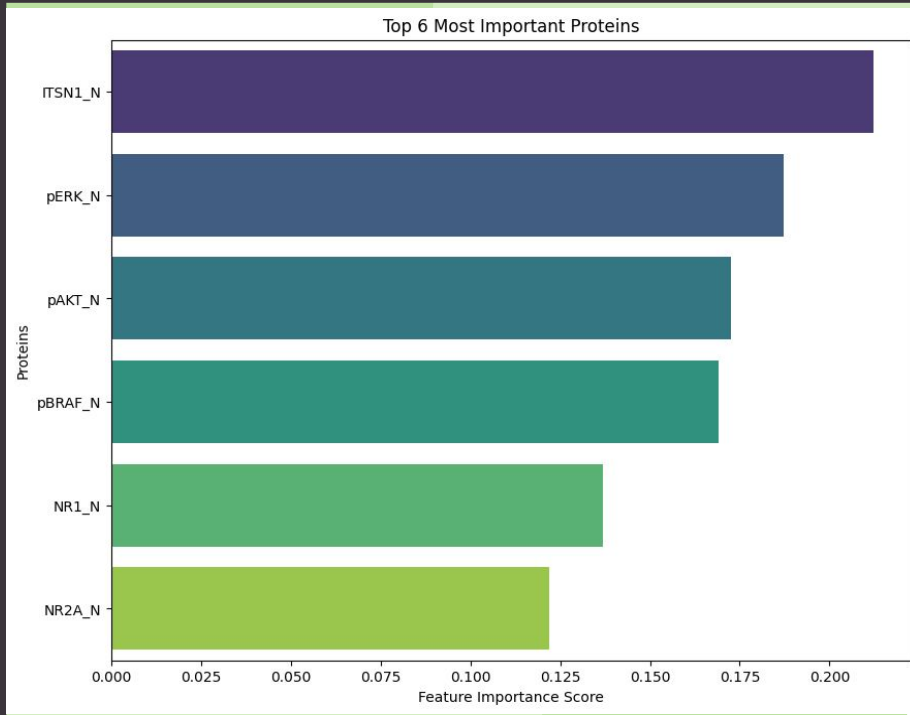
```
Classification Report for Reduced Random Forest Model:
              precision    recall  f1-score   support

     0           0.95         1.00         0.98         21
     1           0.96         0.87         0.92         31
     2           0.97         0.97         0.97         32
     3           0.81         0.96         0.88         27
     4           0.96         0.88         0.92         25
     5           1.00         0.85         0.92         20
     6           0.90         0.88         0.89         32
     7           0.90         1.00         0.95         28

 accuracy          0.93
 macro avg         0.93         0.93         0.93
weighted avg         0.93         0.93         0.93
```

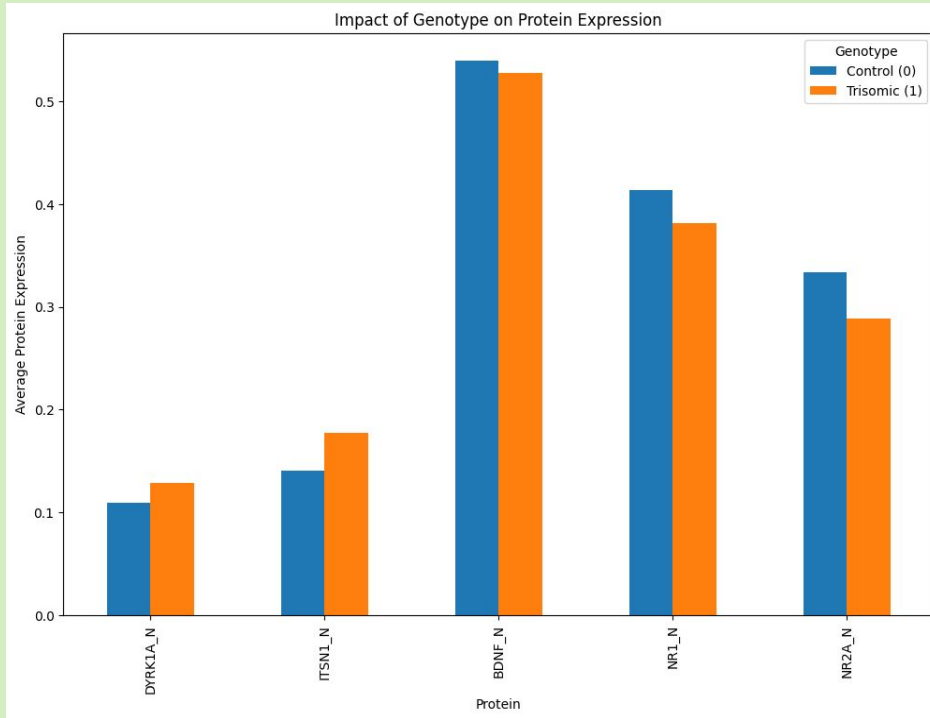
Confusion Matrix for Random Forest





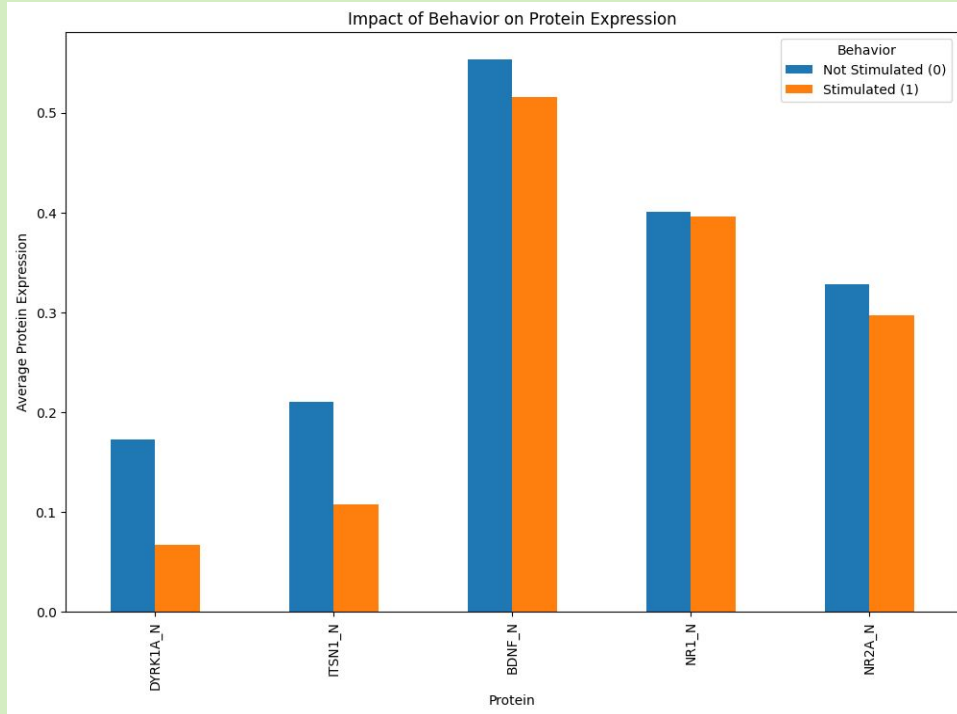
**Top 6 Most
Important
Proteins in Mice
Classification**

Interpretation and Analysis



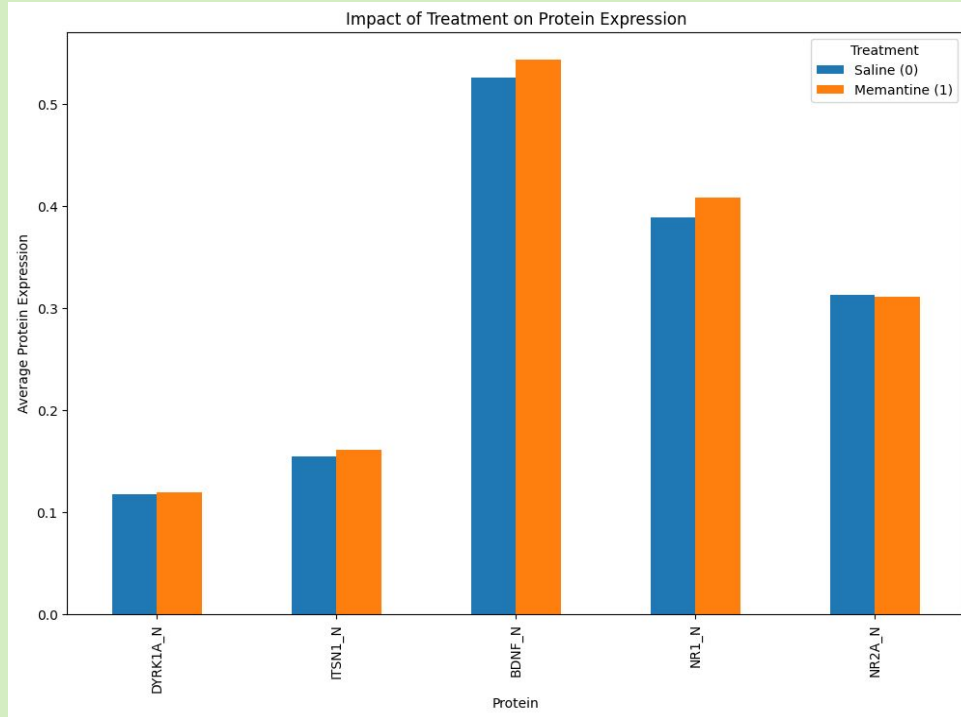
Trisomic mice show higher levels of ITSN1_N and DYRK1A_N, while control mice have higher BDNF_N, NR1_N, and NR2A_N. These protein differences may be linked to the cognitive impairments in Down syndrome.

Interpretation and Analysis



Unstimulated mice show higher levels of DYRK1A_N, ITSN1_N, BDNF_N, and NR2A_N compared to stimulated mice. This suggests that stimulation affects these proteins, potentially influencing associative learning mechanisms.

Interpretation and Analysis



Memantine treatment does not drastically alter most protein expression levels compared to saline. However, slight increases in BDNF_N and NR1_N in the memantine group suggest subtle effects on these proteins, possibly related to its therapeutic effects.

Challenges we faced....




Overfitting

High-dimensional data risks overfitting; regularization and cross-validation are essential.

Missing Data

Technical issues may lead to missing protein values, requiring effective handling techniques.



Biological Variability

Natural variations in protein expression among mice can complicate the identification of subtle class differences.

Conclusion

The study identified key proteins, such as DYRK1A_N, ITSN1_N, and BDNF_N, that differentiate between control and trisomic mice. However, the findings are limited by sample size and biological complexity, suggesting the need for further validation with larger datasets and additional experiments. Future research should focus on confirming these results in broader cohorts, exploring underlying biological mechanisms, and integrating other omics data for a more comprehensive understanding of Down syndrome.



Thank you!

**Does anyone have
any questions?**